



# **Introduction to Meta-Analysis**



# **Introduction to Meta-Analysis**

**Second Edition**

**Michael Borenstein**

*Biosstat, Inc, New Jersey, USA.*

**Larry V. Hedges**

*Northwestern University, Evanston, USA.*

**Julian P.T. Higgins**

*University of Bristol, Bristol, UK.*

**Hannah R. Rothstein**

*Baruch College, New York, USA.*

**WILEY**

This edition first published 2021

© 2021 John Wiley & Sons Ltd

All rights reserved. No part of this publication may be reproduced, stored in a retrieval system, or transmitted, in any form or by any means, electronic, mechanical, photocopying, recording or otherwise, except as permitted by law. Advice on how to obtain permission to reuse material from this title is available at <http://www.wiley.com/go/permissions>.

The right of Michael Borenstein, Larry V. Hedges, Julian P.T. Higgins and Hannah R. Rothstein to be identified as the author of this work has been asserted in accordance with law.

*Registered Offices*

John Wiley & Sons, Inc., 111 River Street, Hoboken, NJ 07030, USA

John Wiley & Sons Ltd, The Atrium, Southern Gate, Chichester, West Sussex, PO19 8SQ, UK

*Editorial Office*

9600 Garsington Road, Oxford, OX4 2DQ, UK

For details of our global editorial offices, customer services, and more information about Wiley products visit us at [www.wiley.com](http://www.wiley.com).

Wiley also publishes its books in a variety of electronic formats and by print-on-demand. Some content that appears in standard print versions of this book may not be available in other formats.

**Limit of Liability/Disclaimer of Warranty**

In view of ongoing research, equipment modifications, changes in governmental regulations, and the constant flow of information relating to the use of experimental reagents, equipment, and devices, the reader is urged to review and evaluate the information provided in the package insert or instructions for each chemical, piece of equipment, reagent, or device for, among other things, any changes in the instructions or indication of usage and for added warnings and precautions. While the publisher and author have used their best efforts in preparing this work, they make no representations or warranties with respect to the accuracy or completeness of the contents of this work and specifically disclaim all warranties, including without limitation any implied warranties of merchantability or fitness for a particular purpose. No warranty may be created or extended by sales representatives, written sales materials or promotional statements for this work. The fact that an organization, website, or product is referred to in this work as a citation and/or potential source of further information does not mean that the publisher and author endorse the information or services the organization, website, or product may provide or recommendations it may make. This work is sold with the understanding that the publisher is not engaged in rendering professional services. The advice and strategies contained herein may not be suitable for your situation. You should consult with a specialist where appropriate. Further, readers should be aware that websites listed in this work may have changed or disappeared between when this work was written and when it is read. Neither the publisher nor author shall be liable for any loss of profit or any other commercial damages, including but not limited to special, incidental, consequential, or other damages.

*Library of Congress Cataloging-in-Publication data applied for*

HB ISBN: 9781119558354

Cover Design: Wiley

Cover Image: Courtesy of Michael Borenstein

Set in 10.5/13pt TimesLTStd by SPI Global, Chennai, India  
Printed and bound by CPI Group (UK) Ltd, Croydon CR0 4YY

C9781119558354\_310321

---

# Contents

---

List of Tables	xv
List of Figures	xix
Acknowledgements	xxv
Preface	xxvii
Preface to the Second Edition	xxxv
Website	xxxvii

## PART 1: INTRODUCTION

<b>1 HOW A META-ANALYSIS WORKS</b>	<b>3</b>
Introduction	3
Individual studies	3
The summary effect	5
Heterogeneity of effect sizes	6
Summary points	7
<b>2 WHY PERFORM A META-ANALYSIS</b>	<b>9</b>
Introduction	9
The streptokinase meta-analysis	10
Statistical significance	11
Clinical importance of the effect	11
Consistency of effects	12
Summary points	13

## PART 2: EFFECT SIZE AND PRECISION

<b>3 OVERVIEW</b>	<b>17</b>
Treatment effects and effect sizes	17
Parameters and estimates	18
Outline of effect size computations	19
<b>4 EFFECT SIZES BASED ON MEANS</b>	<b>21</b>
Introduction	21
Raw (unstandardized) mean difference $D$	21
Standardized mean difference, $d$ and $g$	25
Response ratios	30
Summary points	31

<b>5</b>	<b>EFFECT SIZES BASED ON BINARY DATA (2 × 2 TABLES)</b>	<b>33</b>
Introduction	33	
Risk ratio	33	
Odds ratio	35	
Risk difference	37	
Choosing an effect size index	38	
Summary points	38	
<b>6</b>	<b>EFFECT SIZES BASED ON CORRELATIONS</b>	<b>39</b>
Introduction	39	
Computing $r$	39	
Other approaches	40	
Summary points	41	
<b>7</b>	<b>CONVERTING AMONG EFFECT SIZES</b>	<b>43</b>
Introduction	43	
Converting from the log odds ratio to $d$	44	
Converting from $d$ to the log odds ratio	45	
Converting from $r$ to $d$	45	
Converting from $d$ to $r$	46	
Summary points	47	
<b>8</b>	<b>FACTORS THAT AFFECT PRECISION</b>	<b>49</b>
Introduction	49	
Factors that affect precision	50	
Sample size	50	
Study design	51	
Summary points	53	
<b>9</b>	<b>CONCLUDING REMARKS</b>	<b>55</b>
<b>PART 3: FIXED-EFFECT VERSUS RANDOM-EFFECTS MODELS</b>		
<b>10</b>	<b>OVERVIEW</b>	<b>59</b>
Introduction	59	
Nomenclature	60	
<b>11</b>	<b>FIXED-EFFECT MODEL</b>	<b>61</b>
Introduction	61	
The true effect size	61	
Impact of sampling error	61	
Performing a fixed-effect meta-analysis	63	
Summary points	64	

<b>12 RANDOM-EFFECTS MODEL</b>	<b>65</b>
Introduction	65
The true effect sizes	65
Impact of sampling error	66
Performing a random-effects meta-analysis	68
Summary points	70
<b>13 FIXED-EFFECT VERSUS RANDOM-EFFECTS MODELS</b>	<b>71</b>
Introduction	71
Definition of a summary effect	71
Estimating the summary effect	72
Extreme effect size in a large study or a small study	73
Confidence interval	73
The null hypothesis	76
Which model should we use?	76
Model should not be based on the test for heterogeneity	78
Concluding remarks	79
Summary points	79
<b>14 WORKED EXAMPLES (PART 1)</b>	<b>81</b>
Introduction	81
Worked example for continuous data (Part 1)	81
Worked example for binary data (Part 1)	85
Worked example for correlational data (Part 1)	90
Summary points	94
<b>PART 4: HETEROGENEITY</b>	
<b>15 OVERVIEW</b>	<b>97</b>
Introduction	97
Nomenclature	98
Worked examples	98
<b>16 IDENTIFYING AND QUANTIFYING HETEROGENEITY</b>	<b>99</b>
Introduction	99
Isolating the variation in true effects	99
Computing $Q$	101
Estimating $\tau^2$	106
The $I^2$ statistic	109
Comparing the measures of heterogeneity	111
Confidence intervals for $\tau^2$	114
Confidence intervals (or uncertainty intervals) for $I^2$	115
Summary points	116

<b>17 PREDICTION INTERVALS</b>	<b>119</b>
Introduction	119
Prediction intervals in primary studies	119
Prediction intervals in meta-analysis	121
Confidence intervals and prediction intervals	123
Comparing the confidence interval with the prediction interval	123
Summary points	125
<b>18 WORKED EXAMPLES (PART 2)</b>	<b>127</b>
Introduction	127
Worked example for continuous data (Part 2)	127
Worked example for binary data (Part 2)	131
Worked example for correlational data (Part 2)	134
Summary points	138
<b>19 AN INTUITIVE LOOK AT HETEROGENEITY</b>	<b>139</b>
Introduction	139
Motivating example	140
The $Q$ -value and the $p$ -value do not tell us how much the effect size varies	141
The confidence interval does not tell us how much the effect size varies	142
The $I^2$ statistic does not tell us how much the effect size varies	142
What $I^2$ tells us	142
The $I^2$ index vs. the prediction interval	145
The prediction interval	145
Prediction interval is clear, concise, and relevant	147
Computing the prediction interval	147
How to use $I^2$	149
How to explain heterogeneity	149
How much does the effect size vary across studies?	150
Caveats	150
Conclusion	150
Further reading	151
Summary points	151
The meaning of $I^2$ in Figure 19.2	151
<b>20 CLASSIFYING HETEROGENEITY AS LOW, MODERATE, OR HIGH</b>	<b>155</b>
Introduction	155
Interest should generally focus on an index of absolute heterogeneity	155
The classifications lead themselves to mistakes of interpretation	158
Classifications focus attention in the wrong direction	158
Summary points	158

**PART 5: EXPLAINING HETEROGENEITY**

<b>21 SUBGROUP ANALYSES</b>	<b>161</b>
Introduction	161
Fixed-effect model within subgroups	163
Computational models	172
Random effects with separate estimates of $\tau^2$	174
Random effects with pooled estimate of $\tau^2$	181
The proportion of variance explained	189
Mixed-effects model	192
Obtaining an overall effect in the presence of subgroups	193
Summary points	195
<b>22 META-REGRESSION</b>	<b>197</b>
Introduction	197
Fixed-effect model	198
Fixed or random effects for unexplained heterogeneity	203
Random-effects model	206
Summary points	212
<b>23 NOTES ON SUBGROUP ANALYSES AND META-REGRESSION</b>	<b>213</b>
Introduction	213
Computational model	213
Multiple comparisons	216
Software	216
Analyses of subgroups and regression analyses are observational	217
Statistical power for subgroup analyses and meta-regression	218
Summary points	219

**PART 6: PUTTING IT ALL IN CONTEXT**

<b>24 LOOKING AT THE WHOLE PICTURE</b>	<b>223</b>
Introduction	223
Methylphenidate for adults with ADHD	226
Impact of GLP-1 mimetics on blood pressure	228
Augmenting clozapine with a second antipsychotic	228
Conclusions	231
Caveats	231
Summary points	232
<b>25 LIMITATIONS OF THE RANDOM-EFFECTS MODEL</b>	<b>233</b>
Introduction	233
Assumptions of the random-effects model	234

A textbook case	234
When studies are pulled from the literature	235
A useful fiction	237
Transparency	238
A narrowly defined universe	238
Two important caveats	239
In context	239
Extreme cases	240
Summary points	241
<b>26 KNAPP–HARTUNG ADJUSTMENT</b>	<b>243</b>
Introduction	243
Adjustment is rarely employed in simple analyses	243
Adjusting the standard error	244
The Knapp–Hartung adjustment for other effect size indices	246
$t$ distribution vs. $Z$ distribution	247
Limitations of the Knapp–Hartung adjustment	248
Summary points	249
<b>PART 7: COMPLEX DATA STRUCTURES</b>	
<b>27 OVERVIEW</b>	<b>253</b>
<b>28 INDEPENDENT SUBGROUPS WITHIN A STUDY</b>	<b>255</b>
Introduction	255
Combining across subgroups	255
Comparing subgroups	260
Summary points	260
<b>29 MULTIPLE OUTCOMES OR TIME-POINTS WITHIN A STUDY</b>	<b>263</b>
Introduction	263
Combining across outcomes or time-points	264
Comparing outcomes or time-points within a study	270
Summary points	275
<b>30 MULTIPLE COMPARISONS WITHIN A STUDY</b>	<b>277</b>
Introduction	277
Combining across multiple comparisons within a study	277
Differences between treatments	278
Summary points	279
<b>31 NOTES ON COMPLEX DATA STRUCTURES</b>	<b>281</b>
Introduction	281
Summary effect	281
Differences in effect	282

**PART 8: OTHER ISSUES**

<b>32 OVERVIEW</b>	<b>287</b>
<b>33 VOTE COUNTING – A NEW NAME FOR AN OLD PROBLEM</b>	<b>289</b>
Introduction	289
Why vote counting is wrong	290
Vote counting is a pervasive problem	291
Summary points	293
<b>34 POWER ANALYSIS FOR META-ANALYSIS</b>	<b>295</b>
Introduction	295
A conceptual approach	295
In context	299
When to use power analysis	300
Planning for precision rather than for power	301
Power analysis in primary studies	301
Power analysis for meta-analysis	304
Power analysis for a test of homogeneity	309
Summary points	312
<b>35 PUBLICATION BIAS</b>	<b>313</b>
Introduction	313
The problem of missing studies	314
Methods for addressing bias	316
Illustrative example	317
The model	317
Getting a sense of the data	318
Is there evidence of any bias?	320
How much of an impact might the bias have?	320
Summary of the findings for the illustrative example	324
Conflating bias with the small-study effect	325
Using logic to disentangle bias from small-study effects	326
These methods do not give us the ‘correct’ effect size	327
Some important caveats	327
Procedures do not apply to studies of prevalence	328
The model for publication bias is simplistic	328
Concluding remarks	329
Putting it all together	330
Summary points	330
<b>PART 9: ISSUES RELATED TO EFFECT SIZE</b>	
<b>36 OVERVIEW</b>	<b>335</b>

<b>37 EFFECT SIZES RATHER THAN <i>p</i>-VALUES</b>	<b>337</b>
Introduction	337
Relationship between <i>p</i> -values and effect sizes	337
The distinction is important	339
The <i>p</i> -value is often misinterpreted	340
Narrative reviews vs. meta-analyses	341
Summary points	342
<b>38 SIMPSON'S PARADOX</b>	<b>343</b>
Introduction	343
Circumcision and risk of HIV infection	343
An example of the paradox	345
Summary points	348
<b>39 GENERALITY OF THE BASIC INVERSE-VARIANCE METHOD</b>	<b>349</b>
Introduction	349
Other effect sizes	350
Other methods for estimating effect sizes	353
Individual participant data meta-analyses	354
Bayesian approaches	355
Summary points	357
<b>PART 10: FURTHER METHODS</b>	
<b>40 OVERVIEW</b>	<b>361</b>
<b>41 META-ANALYSIS METHODS BASED ON DIRECTION AND <i>p</i>-VALUES</b>	<b>363</b>
Introduction	363
Vote counting	363
The sign test	363
Combining <i>p</i> -values	364
Summary points	368
<b>42 FURTHER METHODS FOR DICHOTOMOUS DATA</b>	<b>369</b>
Introduction	369
Mantel–Haenszel method	369
One-step (Peto) formula for odds ratio	373
Summary points	376
<b>43 PSYCHOMETRIC META-ANALYSIS</b>	<b>377</b>
Introduction	377
The attenuating effects of artifacts	378
Meta-analysis methods	380
Example of psychometric meta-analysis	381
Comparison of artifact correction with meta-regression	384

---

Sources of information about artifact values	384
How heterogeneity is assessed	385
Reporting in psychometric meta-analysis	386
Concluding remarks	386
Summary points	387
<b>PART 11: META-ANALYSIS IN CONTEXT</b>	
<b>44 OVERVIEW</b>	<b>391</b>
<b>45 WHEN DOES IT MAKE SENSE TO PERFORM A META-ANALYSIS?</b>	<b>393</b>
Introduction	393
Are the studies similar enough to combine?	394
Can I combine studies with different designs?	395
How many studies are enough to carry out a meta-analysis?	399
Summary points	400
<b>46 REPORTING THE RESULTS OF A META-ANALYSIS</b>	<b>401</b>
Introduction	401
The computational model	402
Forest plots	402
Sensitivity analysis	404
Summary points	405
<b>47 CUMULATIVE META-ANALYSIS</b>	<b>407</b>
Introduction	407
Why perform a cumulative meta-analysis?	409
Summary points	412
<b>48 CRITICISMS OF META-ANALYSIS</b>	<b>413</b>
Introduction	413
One number cannot summarize a research field	414
The file drawer problem invalidates meta-analysis	414
Mixing apples and oranges	415
Garbage in, garbage out	416
Important studies are ignored	417
Meta-analysis can disagree with randomized trials	417
Meta-analyses are performed poorly	420
Is a narrative review better?	420
Concluding remarks	422
Summary points	422
<b>49 COMPREHENSIVE META-ANALYSIS SOFTWARE</b>	<b>425</b>
Introduction	425
Features in CMA	426

Teaching elements	427
Documentation	427
Availability	427
Acknowledgments	427
Motivating example	428
Data entry	428
Basic analysis	429
What is the <i>average</i> effect size?	430
How much does the effect size vary?	430
Plot showing distribution of effects	431
High-resolution plot	432
Subgroup analysis	433
Meta-regression	435
Publication bias	438
Explaining results	439
<b>50 HOW TO EXPLAIN THE RESULTS OF AN ANALYSIS</b>	<b>443</b>
Introduction	443
The overview	444
The mean effect size	444
Variation in effect size	444
Notations	444
Impact of resistance exercise on pain	445
Correlation between letter knowledge and word recognition	450
Statins for prevention of cardiovascular events	455
Bupropion for smoking cessation	460
Mortality following mitral-valve procedures in elderly patients	465
<b>PART 12: RESOURCES</b>	
<b>51 SOFTWARE FOR META-ANALYSIS</b>	<b>471</b>
Comprehensive meta-analysis	471
Metafor	471
Stata	472
Revman	472
<b>52 WEB SITES, SOCIETIES, JOURNALS, AND BOOKS</b>	<b>473</b>
Web sites	473
Professional societies	476
Journals	476
Special issues dedicated to meta-analysis	477
Books on systematic review methods and meta-analysis	477
<b>REFERENCES</b>	<b>479</b>
<b>INDEX</b>	<b>491</b>

---

# List of Tables

---

Table 3.1	Roadmap of formulas in subsequent chapters	19
Table 5.1	Nomenclature for $2 \times 2$ table of outcome by treatment	34
Table 5.2	Fictional data for a $2 \times 2$ table	34
Table 8.1	Impact of sample size on variance	51
Table 8.2	Impact of study design on variance	52
Table 14.1	Dataset 1 – Part A (basic data)	82
Table 14.2	Dataset 1 – Part B (fixed-effect computations)	83
Table 14.3	Dataset 1 – Part C (random-effects computations)	85
Table 14.4	Dataset 2 – Part A (basic data)	86
Table 14.5	Dataset 2 – Part B (fixed-effect computations)	87
Table 14.6	Dataset 2 – Part C (random-effects computations)	89
Table 14.7	Dataset 3 – Part A (basic data)	90
Table 14.8	Dataset 3 – Part B (fixed-effect computations)	91
Table 14.9	Dataset 3 – Part C (random-effects computations)	93
Table 16.1	Factors affecting measures of dispersion	111
Table 18.1	Dataset 1 – Part D (intermediate computations)	128
Table 18.2	Dataset 1 – Part E (variance computations)	128
Table 18.3	Dataset 2 – Part D (intermediate computations)	131
Table 18.4	Dataset 2 – Part E (variance computations)	131
Table 18.5	Dataset 3 – Part D (intermediate computations)	135
Table 18.6	Dataset 3 – Part E (variance computations)	135
Table 19.1	Relationship between observed effects and true effects in Figure 19.2, Panel A	152
Table 21.1	Fixed effect model – computations	164
Table 21.2	Fixed-effect model – summary statistics	167
Table 21.3	Fixed-effect model – ANOVA table	169
Table 21.4	Fixed-effect model – subgroups as studies	170
Table 21.5	Random-effects model (separate estimates of $\tau^2$ ) – computations	176
Table 21.6	Random-effects model (separate estimates of $\tau^2$ ) – summary statistics	177
Table 21.7	Random-effects model (separate estimates of $\tau^2$ ) – ANOVA table	180
Table 21.8	Random-effects model (separate estimates of $\tau^2$ ) – subgroups as studies	181

Table 21.9	Statistics for computing a pooled estimate of $\tau^2$	183
Table 21.10	Random-effects model (pooled estimate of $\tau^2$ ) – computations	183
Table 21.11	Random-effects model (pooled estimate of $\tau^2$ ) – summary statistics	185
Table 21.12	Random-effects model (pooled estimate of $\tau^2$ ) – ANOVA table	187
Table 21.13	Random-effects model (pooled estimate of $\tau^2$ ) – subgroups as studies	188
Table 22.1	The BCG dataset	200
Table 22.2	Fixed-effect model – Regression results for BCG	200
Table 22.3	Fixed-effect model – ANOVA table for BCG regression	200
Table 22.4	Random-effects model – regression results for BCG	207
Table 22.5	Random-effects model – test of the model	207
Table 22.6	Random-effects model – comparison of model (latitude) versus the null model	211
Table 26.1	Knapp–Hartung computations for ADHD analysis	244
Table 26.2	Original vs. Knapp–Hartung	246
Table 26.3	Impact of using $t$ distribution on the confidence interval width	248
Table 28.1	Independent subgroups – five fictional studies	256
Table 28.2	Independent subgroups – summary effect	257
Table 28.3	Independent subgroups – synthetic effect for study 1	257
Table 28.4	Independent subgroups – summary effect across studies	258
Table 29.1	Multiple outcomes – five fictional studies	264
Table 29.2	Creating a synthetic variable as the mean of two outcomes	265
Table 29.3	Multiple outcomes – summary effect	267
Table 29.4	Multiple outcomes – impact of correlation on variance of summary effect	269
Table 29.5	Creating a synthetic variable as the difference between two outcomes	271
Table 29.6	Multiple outcomes – difference between outcomes	272
Table 29.7	Multiple outcomes – Impact of correlation on the variance of difference	274
Table 38.1	HIV as function of circumcision (by subgroup)	344
Table 38.2	HIV as function of circumcision – by study	345
Table 38.3	HIV as a function of circumcision – full population	346
Table 38.4	HIV as a function of circumcision – by risk group	346
Table 38.5	HIV as a function of circumcision/risk group – full population	347
Table 39.1	Simple example of a genetic association study	352
Table 41.1	Streptokinase data – calculations for meta-analyses of $p$ -values	367
Table 42.1	Nomenclature for $2 \times 2$ table of events by treatment	370
Table 42.2	Mantel–Haenszel – odds ratio	371

Table 42.3	Mantel–Haenszel – variance of summary effect	372
Table 42.4	One-step – odds ratio and variance	375
Table 43.1	Fictional data for psychometric meta-analysis	382
Table 43.2	Observed (attenuated) correlations	382
Table 43.3	Unattenuated correlations	383



---

# List of Figures

---

Figure 1.1	High-dose versus standard-dose of statins (adapted from Cannon <i>et al.</i> , 2006)	4
Figure 2.1	Impact of streptokinase on mortality (adapted from Lau <i>et al.</i> , 1992)	10
Figure 4.1	Response ratios are analyzed in log units	30
Figure 5.1	Risk ratios are analyzed in log units	34
Figure 5.2	Odds ratios are analyzed in log units	36
Figure 6.1	Correlations are analyzed in Fisher's $z$ units	40
Figure 7.1	Converting among effect sizes	44
Figure 8.1	Impact of sample size on variance	51
Figure 8.2	Impact of study design on variance	52
Figure 10.1	Symbols for true and observed effects	60
Figure 11.1	Fixed-effect model – true effects	62
Figure 11.2	Fixed-effect model – true effects and sampling error	62
Figure 11.3	Fixed-effect model – distribution of sampling error	63
Figure 12.1	Random-effects model – distribution of the true effects	66
Figure 12.2	Random-effects model – true effects	66
Figure 12.3	Random-effects model – true and observed effect in one study	67
Figure 12.4	Random-effects model – between-study and within-study variance	68
Figure 13.1	Fixed-effect model – forest plot showing relative weights	72
Figure 13.2	Random-effects model – forest plot showing relative weights	72
Figure 13.3	Very large studies under fixed-effect model	74
Figure 13.4	Very large studies under random-effects model	74
Figure 14.1	Forest plot of Dataset 1 – fixed-effect weights	84
Figure 14.2	Forest plot of Dataset 1 – random-effects weights	85
Figure 14.3	Forest plot of Dataset 2 – fixed-effect weights	88
Figure 14.4	Forest plot of Dataset 2 – random-effects weights	90
Figure 14.5	Forest plot of Dataset 3 – fixed-effect weights	92
Figure 14.6	Forest plot of Dataset 3 – random-effects weights	94
Figure 16.1	Dispersion across studies relative to error within studies	100
Figure 16.2	$Q$ in relation to $df$ as measure of dispersion	102
Figure 16.3	Flowchart showing how $T^2$ and $I^2$ are derived from $Q$ and $df$	104
Figure 16.4	Impact of $Q$ and number of studies on the $p$ -value	105
Figure 16.5	Impact of excess dispersion and absolute dispersion on $T^2$	107
Figure 16.6	Impact of excess and absolute dispersion on $T$	108

Figure 16.7	Impact of excess dispersion on $I^2$	110
Figure 16.8	Factors affecting $T^2$ but not $I^2$	112
Figure 16.9	Factors affecting $I^2$ but not $T^2$	112
Figure 17.1	Prediction interval based on population parameters $\mu$ and $\tau^2$	121
Figure 17.2	Prediction interval based on sample estimates $M^*$ and $T^2$	122
Figure 17.3	Simultaneous display of confidence interval and prediction interval	123
Figure 17.4	Impact of number of studies on confidence interval and prediction interval	124
Figure 18.1	Forest plot of Dataset 1 – random-effects weights with prediction interval	130
Figure 18.2	Forest plot of Dataset 2 – random-effects weights with prediction interval	134
Figure 18.3	Forest plot of Dataset 3 – random-effects weights with prediction interval	137
Figure 19.1	Alcohol use and mortality. Risk ratio < 1 favors drinkers. Three possible distributions of true effects	141
Figure 19.2	Alcohol use and mortality. Risk ratio < 1 favors drinkers. Three possible distributions of true effects (inner) and observed effects (outer)	144
Figure 19.3	Alcohol use and mortality (Forest plot). Risk ratio < 1 favors drinkers.	146
Figure 19.4	Alcohol use and mortality (true effects). Risk ratio < 1 favors drinkers.	148
Figure 20.1	True effects for two meta-analyses	156
Figure 20.2	True effects (inner) and observed effects (outer) for two meta-analyses	157
Figure 21.1	Fixed-effect model – studies and subgroup effects	164
Figure 21.2	Fixed-effect – subgroup effects	167
Figure 21.3	Fixed-effect model – treating subgroups as studies	170
Figure 21.4	Flowchart for selecting a computational model	174
Figure 21.5	Random-effects model (separate estimates of $\tau^2$ ) – studies and subgroup effects	175
Figure 21.6	Random-effects model (separate estimates of $\tau^2$ ) – subgroup effects	178
Figure 21.7	Random-effects model (separate estimates of $\tau^2$ ) – treating subgroups as studies	180
Figure 21.8	Random-effects model (pooled estimate of $\tau^2$ ) – studies and subgroup effects	182
Figure 21.9	Random-effects model (pooled estimate of $\tau^2$ ) – subgroup effects	185
Figure 21.10	Random-effects model (pooled estimate of $\tau^2$ ) – treating subgroups as studies	188
Figure 21.11	A primary study showing subjects within groups	190

---

Figure 21.12	Random-effects model – variance within and between subgroups	191
Figure 21.13	Proportion of variance explained by subgroup membership	191
Figure 22.1	Fixed-effect model – forest plot for the BCG data	199
Figure 22.2	Fixed-effect model – regression of log risk ratio on latitude	202
Figure 22.3	Fixed-effect model – population effects as function of covariate	204
Figure 22.4	Random-effects model – population effects as a function of covariate	204
Figure 22.5	Random-effects model – forest plot for the BCG data	206
Figure 22.6	Random-effects model – regression of log risk ratio on latitude	208
Figure 22.7	Between-studies variance ( $T^2$ ) with no covariate	210
Figure 22.8	Between-studies variance ( $T^2$ ) with covariate	210
Figure 22.9	Proportion of variance explained by latitude	212
Figure 24.1	Three fictional examples where the mean effect is 0.00	224
Figure 24.2	Three fictional examples where the mean effect is 0.40	225
Figure 24.3	Three fictional examples where the mean effect is 0.80	226
Figure 24.4	Methylphenidate for adults with ADHD (Forest plot). Effect size > 0 favors treatment	227
Figure 24.5	Methylphenidate for adults with ADHD (True effects). Effect size > 0 favors treatment	228
Figure 24.6	GLP-1 mimetics and diastolic BP (Forest plot). Mean difference < 0 favors treatment	229
Figure 24.7	GLP-1 mimetics and diastolic BP (True effects). Mean difference < 0 favors treatment	229
Figure 24.8	Augmenting clozapine (Forest plot). Std mean difference < 0 favors augmentation	230
Figure 24.9	Augmenting clozapine (True effects). Std mean difference < 0 favors augmentation	230
Figure 25.1	Random effects. Confidence interval 60 points wide	234
Figure 25.2	Methylphenidate for adults with ADHD. Effect size > 0 favors treatment	236
Figure 28.1	Creating a synthetic variable from independent subgroups	257
Figure 33.1	The $p$ -value for each study is > 0.20 but the $p$ -value for the summary effect is < 0.02	290
Figure 34.1	Power for a primary study as a function of $n$ and $\delta$	304
Figure 34.2	Power for a meta-analysis as a function of number studies and $\delta$	307
Figure 34.3	Power for a meta-analysis as a function of number studies and heterogeneity	309
Figure 35.1	Passive smoking and lung cancer – forest plot	318
Figure 35.2	Passive smoking and lung cancer – funnel plot	319
Figure 35.3	Observed studies only	321

Figure 35.4	Observed studies and studies imputed by Trim and Fill	322
Figure 35.5	Passive smoking and lung cancer – cumulative forest plot	323
Figure 37.1	Estimating the effect size versus testing the null hypothesis	338
Figure 37.2	The <i>p</i> -value is a poor surrogate for effect size	339
Figure 37.3	Studies where <i>p</i> -values differ but effect sizes is the same	340
Figure 37.4	Studies where <i>p</i> -values are the same but effect sizes differ	341
Figure 37.5	Studies where the more significant <i>p</i> -value corresponds to weaker effect size	341
Figure 38.1	Circumcision and HIV. Odds Ratio > 1 indicates circumcision is associated with lower risk of HIV.	344
Figure 38.2	HIV as function of circumcision – in three sets of studies	347
Figure 41.1	Effect size in four fictional studies	366
Figure 46.1	Forest plot using lines to represent the effect size	403
Figure 46.2	Forest plot using boxes to represent the effect size and relative weight	404
Figure 47.1	Impact of streptokinase on mortality – forest plot	408
Figure 47.2	Impact of streptokinase on mortality – cumulative forest plot	409
Figure 48.1	Forest plot of five fictional studies and a new trial (consistent effects)	418
Figure 48.2	Forest plot of five fictional studies and a new trial (heterogeneous effects)	419
Figure 49.1	Data-entry screen in CMA.	428
Figure 49.2	Basic analysis screen in CMA	429
Figure 49.3	Average effect size (top), Variation in effect size (bottom)	430
Figure 49.4	Plotting distribution of true effects. ADHD	432
Figure 49.5	High-resolution plot in CMA	433
Figure 49.6	Impact of treatment as a function of subgroup: Forest plot	434
Figure 49.7	Impact of treatment as a function of subgroup: Statistics	434
Figure 49.8	Results for regression, random effects	436
Figure 49.9	Regression of effect size on Dose, with SUD held constant	437
Figure 49.10	Funnel plot of observed effects	438
Figure 49.11	Funnel plot of observed and imputed effects	439
Figure 49.12	Regression of effect size ( <i>d</i> ) on Dose and SUD. Plot created in Excel (TM)	441
Figure 50.1	Impact of resistance exercise on pain. Data-entry screen	447
Figure 50.2	Impact of resistance exercise on pain. <i>g</i> > 0 indicates exercise reduced pain	448
Figure 50.3	Impact of resistance exercise on pain. Heterogeneity statistics	449
Figure 50.4	Impact of resistance exercise on pain. Distribution of true effects	449
Figure 50.5	Predicting reading scores. Data-entry screen	452
Figure 50.6	Predicting reading scores	453
Figure 50.7	Predicting reading scores. Heterogeneity statistics	454
Figure 50.8	Predicting reading scores. Distribution of true correlations	454

Figure 50.9	Statins for prevention of cardiovascular events. Data-entry screen	457
Figure 50.10	Statins for prevention of cardiovascular events. Odds ratio < 1 shows reduction in events	458
Figure 50.11	Statins for prevention of cardiovascular events. Heterogeneity statistics	459
Figure 50.12	Statins for prevention of cardiovascular events. Distribution of true effects	459
Figure 50.13	Bupropion for smoking cessation. Data-entry screen	462
Figure 50.14	Bupropion for smoking cessation. Risk ratio > 1 shows reduction in smoking	463
Figure 50.15	Bupropion for smoking cessation. Heterogeneity statistics	464
Figure 50.16	Bupropion for smoking cessation. Distribution of true effects	464
Figure 50.17	Mortality following mitral-valve surgery in elderly patients. Data-entry screen	466
Figure 50.18	Mortality following mitral-valve surgery in elderly patients	467
Figure 50.19	Mortality following mitral-valve surgery in elderly patients. Heterogeneity statistics	468
Figure 50.20	Mortality following mitral-valve surgery in elderly patients. Distribution of true risks	468



---

# Acknowledgements

---

This book was funded by the following grants from the National Institutes of Health: *Combining data types in meta-analysis* (AG021360), *Publication bias in meta-analysis* (AG20052), *Software for meta-regression* (AG024771), from the National Institute on Aging, under the direction of Dr. Sidney Stahl; and *Forest plots for meta-analysis* (DA019280), from the National Institute on Drug Abuse, under the direction of Dr. Thomas Hilton.

These grants allowed us to convene a series of workshops on meta-analysis, and parts of this volume reflect ideas developed as part of these workshops. We would like to acknowledge and thank Doug Altman, Betsy Becker, Jesse Berlin, Michael Brannick, Harris Cooper, Kay Dickersin, Sue Duval, Roger Harbord, Despina Contopoulos-Ioannidis, John Ioannidis, Spyros Konstantopoulos, Mark Lipsey, Mike McDaniel, Ingram Olkin, Fred Oswald, Terri Pigott, Simcha Pollack, David Rindskopf, Stephen Senn, Will Shadish, Jonathan Sterne, Alex Sutton, Thomas Trikalinos, Jeff Valentine, Jack Vevea, Vish Viswesvaran, and David Wilson.

Steven Tarlow helped to edit this book and to ensure the accuracy of all formulas and examples. We would like to acknowledge and thank our editors at Wiley, including Kathryn Sharples, Ashley Alliano, Alison Oliver, Sarah Keegan, Kimberly Monroe-Hill and Viktoria Hartl-Vida. We would especially like to thank Adalfin Jayasingh who served as the production editor for this volume. His attention to detail and his patience in working through revisions are very much appreciated.



---

# Preface

---

In his best-selling book *Baby and Child Care*, Dr. Benjamin Spock wrote 'I think it is preferable to accustom a baby to sleeping on his stomach from the beginning if he is willing'. This statement was included in most editions of the book, and in most of the 50 million copies sold from the 1950s into the 1990s. The advice was not unusual, in that many pediatricians made similar recommendations at the time.

During this same period, from the 1950s into the 1990s, more than 100,000 babies died of sudden infant death syndrome (SIDS), also called *crib death* in the United States and *cot death* in the United Kingdom, where a seemingly healthy baby goes to sleep and never wakes up.

In the early 1990s, researchers became aware that the risk of SIDS decreased by at least 50% when babies were put to sleep on their backs rather than face down. Governments in various countries launched educational initiatives such as the *Back to sleep* campaigns in the United Kingdom and the United States, which led to an immediate and dramatic drop in the number of SIDS deaths.

While the loss of more than 100,000 children would be unspeakably sad in any event, the real tragedy lies in the fact that many of these deaths could have been prevented. Gilbert *et al.* (2005) write

Advice to put infants to sleep on the front for nearly half a century was contrary to evidence available from 1970 that this was likely to be harmful. Systematic review of preventable risk factors for SIDS from 1970 would have led to earlier recognition of the risks of sleeping on the front and might have prevented over 10,000 infant deaths in the UK and at least 50,000 in Europe, the USA and Australasia.

## AN ETHICAL IMPERATIVE

This example is one of several cited by Sir Iain Chalmers in a talk entitled *The scandalous failure of scientists to cumulate scientifically* (Chalmers, 2006). The theme of this talk was that we live in a world where the utility of almost any intervention will be tested repeatedly, and that rather than looking at any study in isolation, we need to look at the body of evidence. While not all systematic reviews carry the urgency of SIDS, the logic of looking at the body of evidence, rather than trying to understand studies in isolation, is always compelling.

Meta-analysis refers to the statistical synthesis of results from a series of studies. While the statistical procedures used in a meta-analysis can be applied to any set of data, the synthesis will be meaningful only if the studies have been collected systematically. This could be in the context of a systematic review, the process of

systematically locating, appraising, and then synthesizing data from a large number of sources. Or, it could be in the context of synthesizing data from a select group of studies, such as those conducted by a pharmaceutical company to assess the efficacy of a new drug.

If a treatment effect (or effect size) is consistent across the series of studies, these procedures enable us to report that the effect is robust across the kinds of populations sampled, and also to estimate the magnitude of the effect more precisely than we could with any of the studies alone. If the treatment effect varies across the series of studies, these procedures enable us to report on the range of effects, and may enable us to identify factors associated with the magnitude of the effect size.

## FROM NARRATIVE REVIEWS TO SYSTEMATIC REVIEWS

Prior to the 1990s, the task of combining data from multiple studies had been primarily the purview of the narrative review. An expert in a given field would read the studies that addressed a question, summarize the findings, and then arrive at a conclusion – for example, that the treatment in question was, or was not, effective. However, this approach suffers from some important limitations.

One limitation is the subjectivity inherent in this approach, coupled with the lack of transparency. For example, different reviewers might use different criteria for deciding which studies to include in the review. Once a set of studies has been selected, one reviewer might give more credence to larger studies, while another gives more credence to ‘quality’ studies and yet another assigns a comparable weight to all studies. One reviewer may require a substantial body of evidence before concluding that a treatment is effective, while another uses a lower threshold. In fact, there are examples in the literature where two narrative reviews come to opposite conclusions, with one reporting that a treatment is effective while the other reports that it is not. As a rule, the narrative reviewer will not articulate (and may not even be fully aware of) the decision-making process used to synthesize the data and arrive at a conclusion.

A second limitation of narrative reviews is that they become *less useful as more information becomes available*. The thought process required for a synthesis requires the reviewer to capture the finding reported in each study, to assign an appropriate weight to that finding, and then to synthesize these findings across all studies in the synthesis. While a reviewer may be able to synthesize data from a few studies in their head, the process becomes difficult and eventually untenable as the number of studies increases. This is true even when the treatment effect (or effect size) is consistent from study to study. Often, however, the treatment effect will vary as a function of study level covariates, such as the patient population, the dose of medication, the outcome variable, and other factors. In these cases, a proper synthesis requires that the researcher be able to understand how the treatment effect varies as a function of these variables, and the narrative review is poorly equipped to address these kinds of issues.

## THE SYSTEMATIC REVIEW AND META-ANALYSIS

For these reasons, beginning in the mid-1980s and taking root in the 1990s, researchers in many fields have been moving away from the narrative review, and adopting systematic reviews and meta-analysis.

For systematic reviews, a clear set of rules is used to search for studies, and then to determine which studies will be included in or excluded from the analysis. Since there is an element of subjectivity in setting these criteria, as well as in the conclusions drawn from the meta-analysis, we cannot say that the systematic review is entirely objective. However, because all of the decisions are specified clearly, the mechanisms are transparent.

A key element in most systematic reviews is the statistical synthesis of the data, or the meta-analysis. Unlike the narrative review, where reviewers implicitly assign some level of importance to each study, in meta-analysis the weights assigned to each study are based on mathematical criteria that are specified in advance. While the reviewers and readers may still differ on the substantive meaning of the results (as they might for a primary study), the statistical analysis provides a transparent, objective, and replicable framework for this discussion.

The formulas used in meta-analysis are extensions of formulas used in primary studies, and are used to address similar kinds of questions to those addressed in primary studies. In primary studies we would typically report a mean and standard deviation for the subjects. If appropriate, we might also use analysis of variance or multiple regression to determine if (and how) subject scores were related to various factors. Similarly, in a meta-analysis, we might report a mean and standard deviation for the treatment effect. And, if appropriate, we would also use procedures analogous to analysis of variance or multiple regression to assess the relationship between the effect and study-level covariates.

Meta-analyses are conducted for a variety of reasons, not only to synthesize evidence on the effects of interventions or to support evidence-based policy or practice. The purpose of the meta-analysis, or more generally, the purpose of any research synthesis, has implications for *when* it should be performed, what model should be used to analyze the data, what sensitivity analyses should be undertaken, and how the results should be interpreted. Losing sight of the fact that meta-analysis is a tool with multiple applications causes confusion and leads to pointless discussions about *what is the right way to perform a research synthesis*, when there is no single right way. It all depends on the purpose of the synthesis, and the data that are available. Much of this book will expand on this idea.

## META-ANALYSIS IS USED IN MANY FIELDS OF RESEARCH

In medicine, systematic reviews and meta-analysis form the core of a movement to ensure that medical treatments are based on the best available empirical data. For example, The Cochrane Collaboration has published the results of over 3700 meta-analyses (as of January 2009) which synthesize data on treatments in all areas

of health care including headaches, cancer, allergies, cardiovascular disease, pain prevention, and depression. The reviews look at interventions relevant to neonatal care, childbirth, infant and childhood diseases, as well as diseases common in adolescents, adults, and the elderly. The kinds of interventions assessed include surgery, drugs, acupuncture, and social interventions. BMJ publishes a series of journals on evidence-based medicine, built on the results from systematic reviews. Systematic reviews and meta-analyses are also used to examine the performance of diagnostic tests, and of epidemiological associations between exposure and disease prevalence, among other topics.

Pharmaceutical companies usually conduct a series of studies to assess the efficacy of a drug. They use meta-analysis to synthesize the data from these studies, yielding a more powerful test (and more precise estimate) of the drug's effect. Additionally, the meta-analysis provides a framework for evaluating the series of studies as a whole, rather than looking at each in isolation. These analyses play a role in internal research, in submissions to governmental agencies, and in marketing. Meta-analyses are also used to synthesize data on adverse events, since these events are typically rare and we need to accumulate information over a series of studies to properly assess the risk of these events.

In the field of education, meta-analysis has been applied to topics as diverse as the comparison of distance education with traditional classroom learning, assessment of the impact of schooling on developing economies, and the relationship between teacher credentials and student achievement. Results of these and similar meta-analyses have influenced practice and policy in various locations around the world.

In psychology, meta-analysis has been applied to basic science as well as in support of evidence-based practice. It has been used to assess personality change over the life span, to assess the influence of media violence on aggressive behavior, and to examine gender differences in mathematics ability, leadership, and nonverbal communication. Meta-analyses of psychological interventions have been used to compare and select treatments for psychological problems, including obsessive-compulsive disorder, impulsivity disorder, bulimia nervosa, depression, phobias, and panic disorder.

In the field of criminology, government agencies have funded meta-analyses to examine the relative effectiveness of various programs in reducing criminal behavior. These include initiatives to prevent delinquency, reduce recidivism, assess the effectiveness of different strategies for police patrols, and for the use of special courts to deal with drug-related crimes.

In business, meta-analyses of the predictive validity of tests that are used as part of the hiring process have led to changes in the types of tests that are used to select employees in many organizations. Meta-analytic results have also been used to guide practices for the reduction of absenteeism, turnover, and counterproductive behavior, and to assess the effectiveness of programs used to train employees.

In the field of ecology, meta-analyses are being used to identify the environmental impact of wind farms, biotic resistance to exotic plant invasion, the effects of changes in the marine food chain, plant reactions to global climate change, the effectiveness of conservation management interventions, and to guide conservation efforts.

## META-ANALYSIS AS PART OF THE RESEARCH PROCESS

Systematic reviews and meta-analyses are used to synthesize the available evidence for a given question to inform policy, as in the examples cited above from medicine, social science, business, ecology, and other fields. While this is probably the most common use of the methodology, meta-analysis can also play an important role in other parts of the research process.

Systematic reviews and meta-analyses can play a role in designing new research. As a first step, they can help determine whether the planned study is necessary. It may be possible to find the required information by synthesizing data from prior studies, and in this case, the research should not be performed. Iain Chalmers (2007) made this point in an article entitled *The lethal consequences of failing to make use of all relevant evidence about the effects of medical treatments: the need for systematic reviews*.

In the event that the new study is needed, the meta-analysis may be useful in helping to design that study. For example, the meta-analysis may show that in the prior studies one outcome index had proven to be more sensitive than others, or that a specific mode of administration had proven to be more effective than others, and should be used in the planned study as well.

For these reasons, various government agencies, including institutes of health in various countries, have been encouraging (or requiring) researchers to conduct a meta-analysis of existing research prior to undertaking new funded studies.

The systematic review can also play a role in the publication of any new primary study. In the introductory section of the publication, a systematic review can help to place the new study in context by describing what we knew before, and what we hoped to learn from the new study. In the discussion section of the publication, a systematic review allows us to address not only the information provided by the new study, but the body of evidence as enhanced by the new study. Iain Chalmers and Michael Clarke (1998) see this approach as a way to avoid studies being reported without context, which they refer to as 'Islands in Search of Continents'. Systematic reviews would provide this context in a more rigorous and transparent manner than the narrative reviews that are typically used for this purpose.

## THE INTENDED AUDIENCE FOR THIS BOOK

Since meta-analysis is a relatively new field, many people, including those who actually use meta-analysis in their work, have not had the opportunity to learn about it systematically. We hope that this volume will provide a framework that allows them to understand the logic of meta-analysis, as well as how to apply and interpret meta-analytic procedures properly.

This book is aimed at researchers, clinicians, and statisticians. Our approach is primarily conceptual. The reader will be able to skip the formulas and still understand, for example, the differences between fixed-effect and random-effects analysis, and the mechanisms used to assess the dispersion in effects from study to study. However, for those with a statistical orientation, we include all the relevant formulas, along with

worked examples. Additionally, the spreadsheets and data files can be downloaded from the web at [www.Introduction-to-Meta-Analysis.com](http://www.Introduction-to-Meta-Analysis.com).

This book can be used as the basis for a course in meta-analysis. Supplementary materials and exercises are posted on the book's website.

This volume is intended for readers from various substantive fields, including medicine, epidemiology, social science, business, ecology, and others. While we have included examples from many of these disciplines, the more important message is that meta-analytic methods that may have developed in any one of these fields have application to all of them.

Since our goal in using these examples is to explain the meta-analysis itself rather than to address the substantive issues, we provide only the information needed for this purpose. For example, we may present an analysis showing that a treatment reduces pain, while ignoring other analyses that show the same treatment increases the risk of adverse events. Therefore, any reader interested in the substantive issues addressed in an example should not rely on this book for that purpose.

#### **AN OUTLINE OF THIS BOOK'S CONTENTS (UPDATED FOR THE SECOND EDITION)**

Part 1 is an introduction to meta-analysis. We present a completed meta-analysis to serve as an example, and highlight the elements of this analysis – the effect size for each study, the summary effect, the dispersion of effects across studies, and so on. Our intent is to show where each element fits into the analysis, and thus provide the reader with a context as they move on to the subsequent parts of the book where each of the elements is explored in detail.

Part 2 introduces the effect sizes, such as the standardized mean difference or the risk ratio, that are computed for each study, and that serve as the unit of currency in the meta-analysis. We also discuss factors that determine the variance of an effect size and show how to compute the variance for each study, since this affects the weight assigned to that study in the meta-analysis.

Part 3 discusses the two computational models used in the vast majority of meta-analyses, the fixed-effect model and the random-effects model. We discuss the conceptual and practical differences between the two, and show how to compute a summary effect using either one.

Part 4 focuses on the issue of dispersion in effect sizes, the fact that the effect size varies from one study to the next. We discuss methods to quantify the heterogeneity, to test it, to incorporate it in the weighting scheme, and to understand it in a substantive as well as a statistical context. In this edition we have expanded this part to address common mistakes in heterogeneity. In particular, we explain that the  $I^2$  statistic is often misinterpreted, and that the practice of classifying dispersion as small, moderate, or high based on  $I^2$  should always be avoided.

Part 5 introduces methods that we might use to understand the reasons for heterogeneity. These include subgroup analyses to compare the effect in different subgroups of studies (analogous to analysis of variance in primary studies), and meta-regression (analogous to multiple regression).

Part 6 is intended to provide context for a meta-analysis. Papers that report a meta-analysis often discuss the mean effect size and heterogeneity as two distinct elements. It is imperative to synthesize the two, and discuss the entire distribution of effects. We discuss the limitations of the random-effects model, and how to take account of these when reporting the results.

Part 7 shows how to work with complex data structures. These include studies that report an effect size for two or more independent subgroups, for two or more outcomes or time-points, and for two or more comparison groups (such as two treatments being compared with the same control).

Part 8 is used to address three separate issues. One chapter discusses the procedure called vote counting, common in narrative reviews, and explains the problems with this approach. One chapter discusses statistical power for a meta-analysis. We show how meta-analysis often (but not always) yields a more powerful test of the null hypothesis than do any of the included studies. Another chapter addresses the question of publication bias. We explain what this is, and discuss methods that have been developed to assess its potential impact.

Part 9 focuses on the issue of why we work with effect sizes in a meta-analysis. In one chapter we explain why we work with effect sizes rather than *p*-values. In another we explain why we compute an effect size for each study, rather than summing data over all studies and then computing an effect size for the summed data. The final chapter in this part shows how the use of inverse-variance weights can be extended to other applications including Bayesian meta-analysis and analyses based on individual participant data.

Part 10 includes chapters on methods that are sometimes used in meta-analysis but that fall outside the central narrative of this volume. These include meta-analyses based on *p*-values, alternate approaches (such as the Mantel–Haenszel method) for assigning study weights, and options sometimes used in psychometric meta-analyses.

Part 11 shows how to take the concepts introduced in Parts 1 to 10 and actually apply them in an analysis. One chapter works through an analysis from start to finish, including a subgroup analysis, meta-regression, and assessment of publication bias. Other chapters present relatively simple analyses using an array of effect-size indices. In all cases we show how to perform the analysis and how to explain the results. We also address the question of when it makes sense to perform a meta-analysis.

Part 12 is a discussion of resources for meta-analysis and systematic reviews. This includes an overview of several computer programs for meta-analysis. It also includes a discussion of organizations that promote the use of systematic reviews and meta-analyses in specific fields, and a list of useful web sites.

## WHAT THIS BOOK DOES NOT COVER

### Other elements of a systematic review

This book deals only with meta-analysis, the statistical formulas and methods used to synthesize data from a set of studies. A meta-analysis can be applied to any data,

but if the goal of the analysis is to provide a synthesis of a body of data from various sources, then it is usually imperative that the data be compiled as part of a systematic review.

A systematic review incorporates many components, such as specification of the question to be addressed, determination of methods to be used for searching the literature and for including or excluding studies, specification of mechanisms to appraise the validity of the included studies, specification of methods to be used for performing the statistical analysis, and a mechanism for disseminating the results.

If the entire review is performed properly, so that the search strategy matches the research question, and yields a reasonably complete and unbiased collection of the relevant studies, then (providing that the included studies are themselves valid) the meta-analysis will also be addressing the intended question. On the other hand, if the search strategy is flawed in concept or execution, or if the studies are providing biased results, then problems exist in the review that the meta-analysis cannot correct.

In Part 12 we include an annotated listing of suggested readings for the other components in the systematic review, but these components are not otherwise addressed in this volume.

### Other meta-analytic methods

In this volume we focus primarily on meta-analyses of effect sizes. That is, analyses where each study yields an estimate of some statistic (a standardized mean difference, a risk ratio, a prevalence, and so on) and our goal is to assess the dispersion in these effects and (if appropriate) compute a summary effect. The vast majority of meta-analyses performed use this approach. We deal only briefly (see Part 10) with other approaches, such as meta-analyses that combine *p*-values rather than effect sizes. We do not address meta-analysis of diagnostic tests, or network meta-analysis.

### Further Reading

- Chalmers, I. (2007). The lethal consequences of failing to make use of all relevant evidence about the effects of medical treatments: the need for systematic reviews. In P. Rothwell (ed.), *Treating Individuals*, ed. London: Lancet: 37–58.
- Chalmers, I., Hedges, L.V. & Cooper, H. (2002). A brief history of research synthesis. *Evaluation in the Health Professions*. 25(1): 12–37.
- Clarke, M., Hopewell, S. & Chalmers, I. (2007). Reports of clinical trials should begin and end with up-to-date systematic reviews of other relevant evidence: a status report. *Journal of the Royal Society of Medicine* 100: 187–190.
- Hunt, M. (1999). *How Science Takes Stock: The Story of Meta-analysis*. New York: Russell Sage Foundation.
- Sutton, A.J. & Higgins, J.P.T. (2008). Recent developments in meta-analysis. *Statistics in Medicine* 27: 625–650.

---

# Preface to the Second Edition

---

The first edition of this text, published in 2009, has been widely embraced by the research community. We are very pleased that this work has informed the practice of meta-analysis, and become a standard text in the field. In this edition we try to improve on that volume in the following ways.

## PRACTICAL INFORMATION

Where the first edition discussed the various statistics that we compute in a meta-analysis, in this edition we show how to use those statistics. More generally, we provide direction for the practical issues that researchers encounter. These issues include the following.

In a meta-analysis to assess the impact of an intervention, the issue of heterogeneity is critically important when we consider the potential utility of the intervention. However, discussions of heterogeneity tend to be superficial. The statistics generally reported for heterogeneity, such as  $Q$ ,  $I^2$ , and  $T^2$ , do not actually tell us how much the effect size varies across studies. Since the reviewer does not have a real understanding of the dispersion in effects, she cannot properly consider the impact of this dispersion. In this edition we show how to quantify dispersion using an intuitive statistic, the prediction interval. This statistic provides information about the dispersion in a clear and concise format. Put simply, the prediction interval provides the information that researchers need, and that they *think* is being provided by other statistics, such as  $I^2$ .

The prediction interval also allows us to consider the mean effect size and the dispersion of effects as a whole, rather than as two separate issues. This allows us to determine, for example, that (a) the intervention is clinically useful in all cases, or (b) the intervention has a substantial benefit in some cases but only a trivial impact in others, or (c) the intervention has a substantial benefit on some cases but is actually harmful in others. Several chapters in this volume provide the foundation for addressing these issues. Additionally, we work through a series of examples to show how to apply these concepts in real analyses.

## LIMITATIONS OF A META-ANALYSIS

In the first edition we explained that when a meta-analysis is based on studies that are pulled from the literature, the random-effects model is usually the one that best fits the analysis. While this is correct, it is important to understand that when we apply the random-effects model for this purpose, there are limitations to what conclusions we can draw. In this edition we discuss those limitations in some detail.

## RECENT DEVELOPMENTS

We have updated the book to keep current with developments in the field of research synthesis.

In the first edition we included a chapter that explains the difference between the two most common statistical models for meta-analysis, the fixed-effect model and the random-effects model. In this edition we have added a discussion of a third model. We have added a chapter on the Knapp–Hartung Sidik–Jonkman adjustment, which applies to confidence intervals and significance tests for the random-effects model. We have also updated the chapter on publication bias.

## HOW TO EXPLAIN THE RESULTS

We have added chapters that provide a “How to” approach to performing and reporting a meta-analysis from start to finish. We have included examples from various fields of research, and using an array of effect-size indices.

Additional information for these examples, including the data sets and step-by-step instructions for performing the analysis using the software *Comprehensive Meta-Analysis* (CMA), is available on the book’s website.

## NEW WEBSITE AND VIDEOS

The book’s website is [www.Introduction-to-Meta-Analysis.com](http://www.Introduction-to-Meta-Analysis.com).

This site includes an array of new features, including videos to illustrate some of the concepts discussed in this volume.

---

# **Website**

---

The book's website is [Introduction-to-Meta-Analysis.com](http://Introduction-to-Meta-Analysis.com)

On this site you can:

- Download the datasets used in this book.
- Watch videos that illustrate some of the concepts in this book.
- Download free software for prediction intervals.
- Download a free trial of Comprehensive Meta-Analysis.
- Submit comments and questions.

For those planning to use this book as a text, there are also worked examples and exercises.

We welcome your feedback.



# **Introduction**



# How a Meta-Analysis Works

---

Introduction

Individual studies

The summary effect

Heterogeneity of effect sizes

---

## INTRODUCTION

Figure 1.1 illustrates a meta-analysis that shows the impact of high dose versus standard dose of statins in preventing death and myocardial infarction (MI). This analysis is adapted from one reported by Cannon *et al.* and published in the *Journal of the American College of Cardiology* (2006).

Our goal in presenting this here is to introduce the various elements in a meta-analysis (the effect size for each study, the weight assigned to each effect size, the estimate of the summary effect, and so on) and show where each fits into the larger scheme. In the chapters that follow, each of these elements will be explored in detail.

## INDIVIDUAL STUDIES

The first four rows on this plot represent the four studies. For each, the study name is shown on the left, followed by the effect size, the relative weight assigned to the study for computing the summary effect, and the *p*-value. The effect size and weight are also shown schematically.

### Effect size

The effect size, a value which reflects the magnitude of the treatment effect or (more generally) the strength of a relationship between two variables, is the unit of currency in a meta-analysis. We compute the effect size for each study, and then work with the effect sizes to assess the consistency of the effect across studies and to compute a summary effect.

### Impact of statin dose on death and myocardial infarction

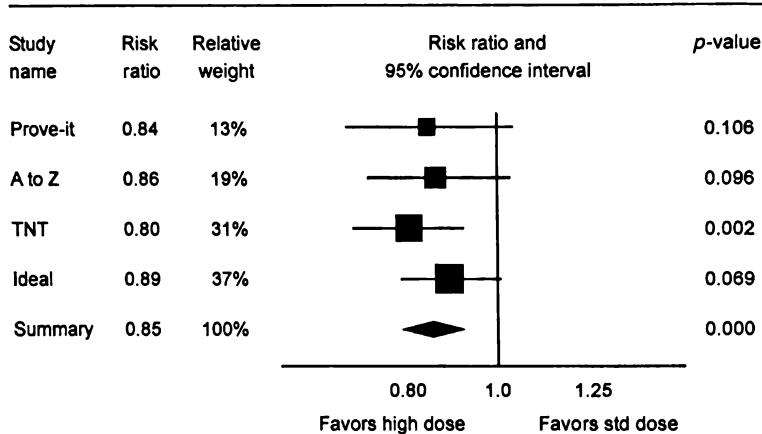


Figure 1.1 High dose versus standard dose of statins (adapted from Cannon *et al.*, 2006).

The effect size could represent the impact of an intervention, such as the impact of medical treatment on risk of infection, the impact of a teaching method on test scores, or the impact of a new protocol on the number of salmon successfully returning upstream. The effect size is not limited to the impact of interventions, but could represent *any relationship* between two variables, such as the difference in test scores for males versus females, the difference in cancer rates for persons exposed or not exposed to second-hand smoke, or the difference in cardiac events for persons with two distinct personality types. In fact, what we generally call an *effect size* could refer simply to the estimate of a single value, such as the prevalence of Lyme disease.

In this example the effect size is the risk ratio. A risk ratio of 1.0 would mean that the risk of death or MI was the same in both groups, while a risk ratio less than 1.0 would mean that the risk was lower in the high-dose group, and a risk ratio greater than 1.0 would mean that the risk was lower in the standard-dose group.

The effect size for each study is represented by a square, with the location of the square representing both the direction and magnitude of the effect. Here, the effect size for each study falls to the left of center (indicating a benefit for the high-dose group). The effect is strongest (most distant from the center) in the *TNT* study and weakest in the *Ideal* study.

Note. For measures of effect size based on ratios (as in this example) a ratio of 1.0 represents no difference between groups. For measures of effect based on differences (such as mean difference), a difference of 0.0 represents no difference between groups.

## Precision

In the schematic, the effect size for each study is bounded by a confidence interval, reflecting the precision with which the effect size has been estimated in that study. The confidence interval for the last study (*Ideal*) is noticeably narrower than that for the first study (*Prove-it*), reflecting the fact that the *Ideal* study has greater precision. The meaning of precision and the factors that affect precision are discussed in Chapter 8.

## Study weights

The solid squares that are used to depict each of the studies vary in size, with the size of each square reflecting the weight that is assigned to the corresponding study when we compute the summary effect. The *TNT* and *Ideal* studies are assigned relatively high weights, while somewhat less weight is assigned to the *A to Z* study and still less to the *Prove-it* study.

As one would expect, there is a relationship between a study's precision and that study's weight in the analysis. Studies with relatively good precision (*TNT* and *Ideal*) are assigned more weight while studies with relatively poor precision (*Prove-it*) are assigned less weight. Since precision is driven primarily by sample size, we can think of the studies as being weighted by sample size.

However, while precision is one of the elements used to assign weights, there are often other elements as well. In Part 3 we discuss different assumptions that one can make about the distribution of effect sizes across studies, and how these affect the weight assigned to each study.

## p-values

For each study we show the *p*-value for a test of the null hypothesis. There is a necessary correspondence between the *p*-value and the confidence interval, such that the *p*-value will fall under 0.05 if and only if the 95% confidence interval does not include the null value. Therefore, by scanning the confidence intervals we can easily identify the statistically significant studies. The role of *p*-values in the analysis, as well as the relationship between *p*-values and effect size, is discussed in Chapter 37.

In this example, for three of the four studies the confidence interval crosses the null hypothesis, and the *p*-value is greater than 0.05. In one (the *TNT* study) the confidence interval does not cross the null hypothesis, and the *p*-value falls under 0.05.

## THE SUMMARY EFFECT

One goal of the synthesis is usually to compute a summary effect. Typically we report the effect size itself, as well as a measure of precision and a *p*-value.

## Effect size

On the plot the summary effect is shown on the bottom line. In this example the summary risk ratio is 0.85, indicating that the risk of death (or MI) was 15% lower for patients assigned to the high dose than for patients assigned to standard dose.

The summary effect is nothing more than the weighted mean of the individual effects. However, the mechanism used to assign the weights (and therefore the meaning of the summary effect) depends on our assumptions about the distribution of effect sizes from which the studies were sampled. Under the fixed-effect model, we assume that all studies in the analysis share the same true effect size, and the summary effect is our estimate of this common effect size. Under the random-effects model, we assume that the true effect size varies from study to study, and the summary effect is our estimate of the mean of the distribution of effect sizes. This is discussed in Part 3.

## Precision

The summary effect is represented by a diamond. The location of the diamond represents the effect size while its width reflects the precision of the estimate. In this example the diamond is centered at 0.85, and extends from 0.79 to 0.92, meaning that the actual impact of the high dose (as compared to the standard) likely falls somewhere in that range.

The precision addresses the accuracy of the summary effect as an estimate of the true effect. However, as discussed in Part 3, the exact meaning of the precision depends on the statistical model.

## *p*-value

The *p*-value for the summary effect is 0.00003. This *p*-value reflects both the magnitude of the summary effect size and also the volume of information on which the estimate is based. Note that the *p*-value for the summary effect is substantially more compelling than that of any single study. Indeed, only one of the four studies had a *p*-value under 0.05. The relationship between *p*-values and effect sizes is discussed in Chapter 37.

## HETEROGENEITY OF EFFECT SIZES

In this example the treatment effect is consistent across all studies (by a criterion explained in Chapter 16), but such is not always the case. A key theme in this volume is the importance of assessing the dispersion of effect sizes from study to study, and then taking this into account when interpreting the data. If the effect size is consistent, then we will usually focus on the summary effect, and note that this effect is robust across the domain of studies included in the analysis. If the effect size varies modestly, then we might still report the summary effect but note that the true effect in any given study could be somewhat lower or higher than this value. If the effect varies

substantially from one study to the next, our attention will shift from the summary effect to the dispersion itself.

Because the dispersion in observed effects is partly spurious (it includes both real difference in effects and also random error), before trying to interpret the variation in effects we need to determine what part (if any) of the observed variation is real. In Part 4 we show how to partition the observed variance into the part due to error and the part that represents variation in true effect sizes, and then how to use this information in various ways.

In this example our goal was to estimate the summary effect in one set of populations. In some cases, however, we will want to compare the effect size for one subgroup of studies versus another (say, for studies that used an elderly population versus those that used a relatively young population). In other cases we may want to assess the impact of putative moderators (or covariates) on the effect size (say, comparing the effect size in studies that used doses of 10, 20, 40, 80, 160 mg.). These kinds of analyses are also discussed in Part 4.

### SUMMARY POINTS

- To perform a meta-analysis we compute an effect size and variance for each study, and then compute a weighted mean of these effect sizes.
- To compute the weighted mean we generally assign more weight to the more precise studies, but the rules for assigning weights depend on our assumptions about the distribution of true effects.



# Why Perform a Meta-Analysis

---

### Introduction

The streptokinase meta-analysis

Statistical significance

Clinical importance of the effect

Consistency of effects

---

## INTRODUCTION

Why perform a meta-analysis? What are the advantages of using statistical methods to synthesize data rather than taking the results that had been reported for each study and then having these collated and synthesized by an expert?

In this chapter we start at the point where we have already selected the studies to be included in the review, and are planning the synthesis itself. We do not address the differences between systematic reviews and narrative reviews in the process of locating and selecting studies. These differences can be critically important, but (as always) our focus is on the data analysis rather than the full process of the review.

The goal of a synthesis is to understand the results of any study in the context of all the other studies. First, we need to know whether or not the effect size is consistent across the body of data. If it *is* consistent, then we want to estimate the effect size as accurately as possible and to report that it is robust across the kinds of studies included in the synthesis. On the other hand, if it varies substantially from study to study, we want to quantify the extent of the variance and consider the implications.

Meta-analysis is able to address these issues whereas the narrative review is not. We start with an example to show how meta-analysis and narrative review would approach the same question, and then use this example to highlight the key differences between the two.

## THE STREPTOKINASE META-ANALYSIS

During the time period beginning in 1959 and ending in 1988 (a span of nearly 30 years) there were a total of 33 randomized trials performed to assess the ability of streptokinase to prevent death following a heart attack. Streptokinase, a so-called *clot buster* which is administered intravenously, was hypothesized to dissolve the clot causing the heart attack, and thus increase the likelihood of survival. The trials all followed similar protocols, with patients assigned at random to either treatment or a placebo. The outcome, whether or not the patient died, was the same in all the studies.

The trials varied substantially in size. The median sample size was slightly over 100 but there was one trial with a sample size in the range of 20 patients, and two large scale trials which enrolled some 12,000 and 17,000 patients, respectively. Of the 33 studies, six were statistically significant while the other 27 were not, leading to the perception that the studies yielded conflicting results.

In 1992 Lau *et al.* published a meta-analysis that synthesized the results from the 33 studies. The presentation that follows is based on the Lau paper (though we use a risk ratio where Lau used an odds ratio).

The forest plot (Figure 2.1) provides context for the analysis. An effect size to the left of center indicates that treated patients were more likely to survive, while an effect size to the right of center indicates that control patients were more likely to survive.

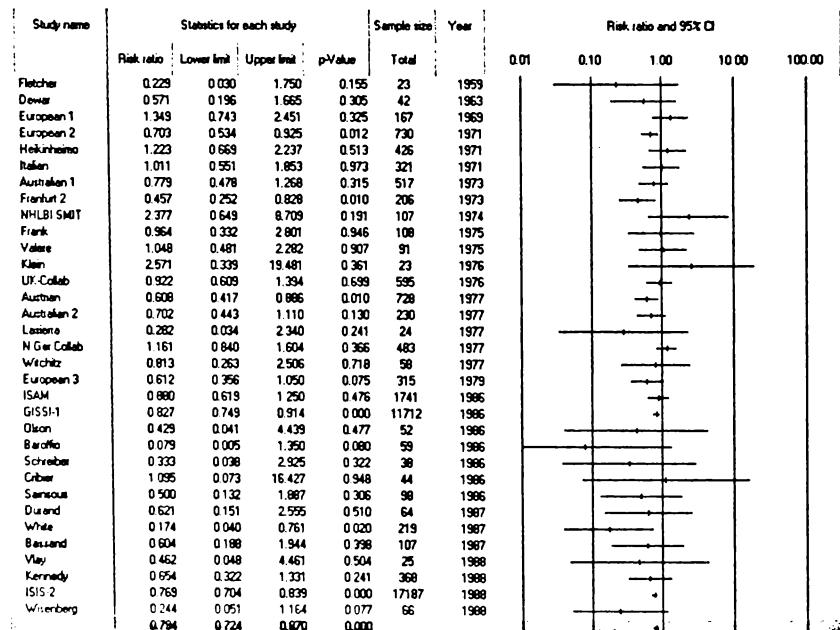


Figure 2.1 Impact of streptokinase on mortality (adapted from Lau *et al.*, 1992).

The plot serves to highlight the following points.

- The effect sizes are reasonably consistent from study to study. Most fall in the range of 0.50 to 0.90, which suggests that it would be appropriate to compute a summary effect size.
- The summary effect is a risk ratio of 0.79 with a 95% confidence interval of 0.72 to 0.87 (that is, a 21% decrease in risk of death, with 95% confidence interval of 13% to 28%). The  $p$ -value for the summary effect is 0.0000008.
- The confidence interval that bounds each effect size indicates the precision in that study. If the interval excludes 1.0, the  $p$ -value is less than 0.05 and the study is statistically significant. Six of the studies were statistically significant while 27 were not.

In sum, the treatment reduces the risk of death by some 21%. And this effect was reasonably consistent across all studies in the analysis.

Over the course of this volume we explain the statistical procedures that led to these conclusions. Our goal in the present chapter is simply to explain that meta-analysis does offer these mechanisms, whereas the narrative review does not. The key differences are as follows.

## STATISTICAL SIGNIFICANCE

One of the first questions asked of a study is the statistical significance of the results. The narrative review has no mechanism for synthesizing the  $p$ -values from the different studies, and must deal with them as discrete pieces of data. In this example six of the studies were statistically significant while the other 27 were not, which led some to conclude that there was evidence against an effect, or that the results were inconsistent (see vote counting in Chapter 33). By contrast, the meta-analysis allows us to combine the effects and evaluate the statistical significance of the summary effect. The  $p$ -value for the summary effect is  $p = 0.0000008$ .

While one might assume that 27 studies failed to reach statistical significance because they reported small effects, it is clear from the forest plot that this is not the case. In fact, the treatment effect in many of these studies was actually *larger* than the treatment effect in the six studies that *were* statistically significant. Rather, the reason that 82% of the studies were not statistically significant is that these studies had small sample sizes and low statistical power. In fact, as discussed in Chapter 34, most had power of less than 20%. By contrast, power for the meta-analysis exceeded 99.9% (see Chapter 34).

As in this example, if the goal of a synthesis is to test the null hypothesis, then meta-analysis provides a mathematically rigorous mechanism for this purpose. However, meta-analysis also allows us to move beyond the question of statistical significance, and address questions that are more interesting and also more relevant.

## CLINICAL IMPORTANCE OF THE EFFECT

Since the point of departure for a narrative review is usually the  $p$ -values reported by the various studies, the review will often focus on the question of whether or not the

body of evidence allows us to reject the null hypothesis. There is no good mechanism for discussing the magnitude of the effect. By contrast, the meta-analytic approaches discussed in this volume allow us to compute an estimate of the effect size for each study, and these effect sizes fall at the core of the analysis.

This is important because the effect size is what we care about. If a clinician or patient needs to make a decision about whether or not to employ a treatment, they want to know if the treatment reduces the risk of death by 5% or 10% or 20%, and this is the information carried by the effect size. Similarly, if we are thinking of implementing an intervention to increase the test scores of students, or to reduce the number of incarcerations among at-risk juveniles, or to increase the survival time for patients with pancreatic cancer, the question we ask is about the magnitude of the effect. The *p*-value can tell us only that the effect is not zero, and to report simply that the effect is not zero is to miss the point.

## CONSISTENCY OF EFFECTS

When we are working with a collection of studies, it is critically important to ask whether or not the effect size is consistent across studies. The implications are quite different for a drug that consistently reduces the risk of death by 20%, as compared with a drug that reduces the risk of death by 20% on average, but that increases the risk by 20% in some populations while reducing it by 60% in others.

The narrative review has no good mechanism for assessing the consistency of effects. The narrative review starts with *p*-values, and because the *p*-value is driven by the size of a study as well as the effect in that study, the fact that one study reported a *p*-value of 0.001 and another reported a *p*-value of 0.50 does not mean that the effect was larger in the former. The *p*-value of 0.001 *could* reflect a large effect size but it could also reflect a moderate or small effect in a large study (see the GISSI-1 study in Figure 2.1, for example). The *p*-value of 0.50 *could* reflect a small (or nil) effect size but could also reflect a large effect in a small study (see the Fletcher study, for example).

This point is often missed in narrative reviews. Often, researchers interpret a non-significant result to mean that there is no effect. If some studies are statistically significant while others are not, the reviewers see the results as conflicting. This problem runs through many fields of research. To borrow a phrase from Cary Grant's character in *Arsenic and Old Lace*, we might say that it practically gallops.

Schmidt (1996) outlines the impact of this practice on research and policy. Suppose an idea is proposed that will improve test scores for African-American children. A number of studies are performed to test the intervention. The effect size is positive and consistent across studies but power is around 50%, and only around 50% of the studies yield statistically significant results. Researchers report that the evidence is 'conflicting' and launch a series of studies to determine why the intervention had a positive effect in some studies but not others (Is it the teacher's attitude? Is it the students' socioeconomic status?), entirely missing the point that the effect was actually

consistent from one study to the next. No pattern can be found (since none exists). Eventually, researchers decide that the issue cannot be understood. A promising idea is lost, and a perception builds that research is not to be trusted. A similar point is made by Meehl (1978, 1990).

Rossi (1997) gives an example from the field of memory research that shows what can happen to a field of research when reviewers work with discrete  $p$ -values. The issue of whether or not researchers could demonstrate the spontaneous recovery of previously extinguished associations had a bearing on a number of important learning theories, and some 40 studies on the topic were published between 1948 and 1969. Evidence of the effect (that is, statistically significant findings) was obtained in only about half the studies, which led most texts and reviews to conclude that the effect was ephemeral and ‘the issue was not so much resolved as it was abandoned’ (p. 179). Later, Rossi returned to these studies and found that the average effect size ( $d$ ) was 0.39. If we assume that this is the population effect size, the mean power for these studies would have been slightly under 50%. On this basis we would expect about half the studies to yield a significant effect, which is exactly what happened.

Even worse, when the significant study was performed in one type of sample and the nonsignificant study was performed in another type of sample, researchers would sometimes interpret this difference as meaning that the effect existed in one population but not the other. Abelson (1997) notes that if a treatment effect yields a  $p$ -value of 0.07 for wombats and 0.05 for dingbats we are likely to see a discussion explaining why the treatment is effective only in the latter group—completely missing the point that the treatment effect may have been virtually identical in the two. The treatment effect may have even been *larger* for the wombats if the sample size was smaller.

By contrast, meta-analysis completely changes the landscape. First, we work with effect sizes (not  $p$ -values) to determine whether or not the effect size is consistent across studies. Additionally, we apply methods based on statistical theory to allow that some (or all) of the observed dispersion is due to random sampling variation rather than differences in the true effect sizes. Then, we apply formulas to partition the variance into random error versus real variance, to quantify the true differences among studies, and to consider the implications of this variance. In the Schmidt and the Rossi examples, a meta-analysis might have found that the effect size was consistent across studies, and that all of the observed variation in effects could be attributed to random sampling error.

### SUMMARY POINTS

- Since the narrative review is based on discrete reports from a series of studies, it provides no real mechanism for synthesizing the data. To borrow a phrase from Abelson, it involves *doing arithmetic with words*. And, when the words are based on  $p$ -values *the words are the wrong words*.

- By contrast, in a meta-analysis we introduce two fundamental changes. First, we work directly with the effect size from each study rather than the  $p$ -value. Second, we include all of the effects in a single statistical synthesis. This is critically important for the goal of computing (and testing) a summary effect. Meta-analysis also allows us to assess the dispersion of effects, and distinguish between real dispersion and spurious dispersion.

# Effect Size and Precision



# Overview

---

Treatment effects and effect sizes

Parameters and estimates

Outline of effect size computations

---

### TREATMENT EFFECTS AND EFFECT SIZES

The terms *treatment effects* and *effect sizes* are used in different ways by different people. Meta-analyses in medicine often refer to the effect size as a *treatment effect*, and this term is sometimes assumed to refer to odds ratios, risk ratios, or risk differences, which are common in meta-analyses that deal with medical interventions. Similarly, meta-analyses in the social sciences often refer to the effect size simply as an *effect size* and this term is sometimes assumed to refer to standardized mean differences or to correlations, which are common in social science meta-analyses.

In fact, though, both the terms *effect size* and *treatment effect* can refer to any of these indices, and the distinction between these terms lies not in the index itself but rather in the nature of the study. The term *effect size* is appropriate when the index is used to quantify the relationship between two variables or a difference between two groups. By contrast, the term *treatment effect* is appropriate only for an index used to quantify the impact of a deliberate intervention. Thus, the difference between males and females could be called an *effect size* only, while the difference between treated and control groups could be called either an *effect size* or a *treatment effect*.

While most meta-analyses focus on relationships between variables, some have the goal of estimating a mean or risk or rate in a single population. For example, a meta-analysis might be used to combine several estimates for the prevalence of Lyme disease in Wabash or the mean SAT score for students in Utah. In these cases the index is clearly not a treatment effect, and is also not an effect size, since *effect* implies a relationship. Rather, the parameter being estimated could be called simply a *single group summary*.

Note, however, that the classification of an index as an *effect size* and/or a *treatment effect* (or simply a *single group summary*) has no bearing on the computations. In the

meta-analysis itself we have simply a series of values and their variances, and the same mathematical formulas apply. In this volume we generally use the term *effect size*, but we use it in a generic sense, to include also treatment effects, single group summaries, or even a generic statistic.

### How to choose an effect-size index

Three major considerations should drive the choice of an effect size index. The first is that the effect sizes from the different studies should be comparable to one another in the sense that they measure (at least approximately) the same thing. That is, the effect size should not depend on aspects of study design that may vary from study to study (such as sample size or whether covariates are used). The second is that estimates of the effect size should be computable from the information that is likely to be reported in published research reports. That is, it should not require the re-analysis of the raw data (unless these are known to be available). The third is that the effect size should have good technical properties. For example, its sampling distribution should be known so that variances and confidence intervals can be computed.

Additionally, the effect size should be substantively interpretable. This means that researchers in the substantive area of the work represented in the synthesis should find the effect size meaningful. If the effect size is not inherently meaningful, it is usually possible to transform the effect size to another metric for presentation. For example, the analyses may be performed using the log risk ratio but then transformed to a risk ratio (or even to Illustrative absolute risks) for presentation.

In practice, the kind of data used in the primary studies will usually lead to a pool of two or three effect sizes that meet the criteria outlined above, which makes the process of selecting an effect size relatively straightforward. If the summary data reported by the primary study are based on means and standard deviations in two groups, the appropriate effect size will usually be either the raw difference in means, the standardized difference in means, or the response ratio. If the summary data are based on a binary outcome such as events and non-events in two groups, the appropriate effect size will usually be the risk ratio, the odds ratio, or the risk difference. If the primary study reports a correlation between two variables, then the correlation coefficient itself may serve as the effect size.

## PARAMETERS AND ESTIMATES

Throughout this volume we make the distinction between an underlying effect size parameter (denoted by the Greek letter  $\theta$ ) and the sample estimate of that parameter (denoted by  $Y$ ).

If a study had an infinitely large sample size then it would yield an effect size  $Y$  that was identical to the population parameter  $\theta$ . In fact, though, sample sizes are finite and so the effect size estimate  $Y$  always differs from  $\theta$  by some amount. The value of  $Y$  will vary from sample to sample, and the distribution of these values is the sampling distribution of  $Y$ . Statistical theory allows us to estimate the sampling distribution of effect size estimates, and hence their standard errors.

## OUTLINE OF EFFECT SIZE COMPUTATIONS

Table 3.1 provides an outline of the computational formulas that follow.

These are some of the more common effect sizes and study designs. A more extensive array of formulas is offered in Borenstein *et al.* (in preparation).

**Table 3.1** Roadmap of formulas in subsequent chapters.

---

**Effect sizes based on means (Chapter 4)**

Raw (unstandardized) mean difference ( $D$ )

    Based on studies with independent groups

    Based on studies with matched groups or pre-post designs

Standardized mean difference ( $d$  or  $g$ )

    Based on studies with independent groups

    Based on studies with matched groups or pre-post designs

Response ratios ( $R$ )

    Based on studies with independent groups

**Effect sizes based on binary data (Chapter 5)**

Risk ratio ( $RR$ )

    Based on studies with independent groups

Odds ratio ( $OR$ )

    Based on studies with independent groups

Risk difference ( $RD$ )

    Based on studies with independent groups

**Effect sizes based on correlational data (Chapter 6)**

Correlation ( $r$ )

    Based on studies with one group

---



# Effect Sizes Based on Means

---

### Introduction

Raw (unstandardized) mean difference  $D$

Standardized mean difference,  $d$  and  $g$

Response ratios

---

## INTRODUCTION

When the studies report means and standard deviations, the preferred effect size is usually the raw mean difference, the standardized mean difference, or the response ratio. These effect sizes are discussed in this chapter.

### RAW (UNSTANDARDIZED) MEAN DIFFERENCE $D$

When the outcome is reported on a meaningful scale *and* all studies in the analysis use the same scale, the meta-analysis can be performed directly on the raw difference in means (henceforth, we will use the more common term, *raw mean difference*). The primary advantage of the raw mean difference is that it is intuitively meaningful, either inherently (for example, blood pressure, which is measured on a known scale) or because of widespread use (for example, a national achievement test for students, where all relevant parties are familiar with the scale).

Consider a study that reports means for two groups (Treated and Control) and suppose we wish to compare the means of these two groups. Let  $\mu_1$  and  $\mu_2$  be the true (population) means of the two groups. The population mean difference is defined as

$$\Delta = \mu_1 - \mu_2. \tag{4.1}$$

In the two sections that follow we show how to compute an estimate  $D$  of this parameter and its variance from studies that used two independent groups and from studies that used paired groups or matched designs.

### Computing $D$ from studies that use independent groups

We can estimate the mean difference  $\Delta$  from a study that used two independent groups as follows. Let  $\bar{X}_1$  and  $\bar{X}_2$  be the sample means of the two independent groups. The sample estimate of  $\Delta$  is just the difference in sample means, namely

$$D = \bar{X}_1 - \bar{X}_2. \quad (4.2)$$

Note that uppercase  $D$  is used for the *raw* mean difference, whereas lowercase  $d$  will be used for the *standardized* mean difference (below).

Let  $S_1$  and  $S_2$  be the sample standard deviations of the two groups, and  $n_1$  and  $n_2$  be the sample sizes in the two groups. If we assume that the two population standard deviations are the same (as is assumed to be the case in most parametric data analysis techniques), so that  $\sigma_1 = \sigma_2 = \sigma$ , then the variance of  $D$  is

$$V_D = \frac{n_1 + n_2}{n_1 n_2} S_{\text{pooled}}^2, \quad (4.3)$$

where

$$S_{\text{pooled}} = \sqrt{\frac{(n_1 - 1)S_1^2 + (n_2 - 1)S_2^2}{n_1 + n_2 - 2}}. \quad (4.4)$$

If we don't assume that the two population standard deviations are the same, then the variance of  $D$  is

$$V_D = \frac{S_1^2}{n_1} + \frac{S_2^2}{n_2}. \quad (4.5)$$

In either case, the standard error of  $D$  is then the square root of  $V$ ,

$$SE_D = \sqrt{V_D}. \quad (4.6)$$

For example, suppose that a study has sample means  $\bar{X}_1 = 103.00$ ,  $\bar{X}_2 = 100.00$ , sample standard deviations  $S_1 = 5.5$ ,  $S_2 = 4.5$ , and sample sizes  $n_1 = n_2 = 50$ . The raw mean difference  $D$  is

$$D = 103.00 - 100.00 = 3.00.$$

If we assume that  $\sigma_1 = \sigma_2$  then the pooled standard deviation within groups is

$$S_{\text{pooled}} = \sqrt{\frac{(50 - 1) \times 5.5^2 + (50 - 1) \times 4.5^2}{50 + 50 - 2}} = 5.0249.$$

The variance and standard error of  $D$  are given by

$$V_D = \frac{50 + 50}{50 \times 50} \times 5.0249^2 = 1.0100,$$

and

$$SE_D = \sqrt{1.0100} = 1.0050.$$

If we do not assume that  $\sigma_1 = \sigma_2$  then the variance and standard error of  $D$  are given by

$$V_D = \frac{5.5^2}{50} + \frac{4.5^2}{50} = 1.0100$$

and

$$SE_D = \sqrt{1.0100} = 1.0050.$$

In this example formulas (4.3) and (4.5) yield the same result, but this will be true only if the sample size and/or the estimate of the variances is the same in the two groups.

### Computing $D$ from studies that use matched groups or pre-post scores

The previous formulas are appropriate for studies that use two independent groups. Another study design is the use of matched groups, where pairs of participants are matched in some way (for example, siblings, or patients at the same stage of disease), with the two members of each pair then being assigned to different groups. The unit of analysis is the pair, and the advantage of this design is that each pair serves as its own control, reducing the error term and increasing the statistical power. The magnitude of the impact depends on the correlation between (for example) siblings, with a higher correlation yielding a lower variance (and increased precision).

The sample estimate of  $\Delta$  is just the sample mean difference,  $D$ . If we have the difference score for each pair, which gives us the mean difference  $\bar{X}_{diff}$  and the standard deviation of these differences ( $S_{diff}$ ), then

$$D = \bar{X}_{diff}, \quad (4.7)$$

$$V_D = \frac{S_{diff}^2}{n}, \quad (4.8)$$

where  $n$  is the number of pairs, and

$$SE_D = \sqrt{V_D}. \quad (4.9)$$

For example, if the mean difference is 5.00 with standard deviation of the difference of 10.00 and  $n$  of 50 pairs, then

$$D = 5.0000,$$

$$V_D = \frac{10.00^2}{50} = 2.0000, \quad (4.10)$$

and

$$SE_D = \sqrt{2.00} = 1.4142. \quad (4.11)$$

Alternatively, if we have the mean and standard deviation for each set of scores (for example, siblings  $A$  and  $B$ ), the difference is

$$D = \bar{X}_1 - \bar{X}_2. \quad (4.12)$$

The variance is again given by

$$V_D = \frac{S_{diff}^2}{n}, \quad (4.13)$$

where  $n$  is the number of pairs, and the standard error is given by

$$SE_D = \sqrt{V_D}. \quad (4.14)$$

However, in this case we need to compute the standard deviation of the difference scores from the standard deviation of each sibling's scores. This is given by

$$S_{diff} = \sqrt{S_1^2 + S_2^2 - 2 \times r \times S_1 \times S_2} \quad (4.15)$$

where  $r$  is the correlation between 'siblings' in matched pairs. If  $S_1 = S_2$ , then (4.15) simplifies to

$$S_{diff} = \sqrt{2 \times S_{pooled}^2(1 - r)}. \quad (4.16)$$

In either case, as  $r$  moves toward 1.0 the standard error of the paired difference will decrease, and when  $r = 0$  the standard error of the difference is the same as it would be for a study with two independent groups, each of size  $n$ .

For example, suppose the means for siblings A and B are 105.00 and 100.00, with standard deviations 10 and 10, the correlation between the two sets of scores is 0.50, and the number of pairs is 50. Then

$$D = 105.00 - 100.00 = 5.0000,$$

$$V_D = \frac{10.00^2}{50} = 2.0000,$$

and

$$SE_D = \sqrt{2.00} = 1.4142.$$

In the calculation of  $V_d$ , the  $S_{diff}$  is computed using

$$S_{diff} = \sqrt{10^2 + 10^2 - 2 \times 0.50 \times 10 \times 10} = 10.0000$$

or

$$S_{diff} = \sqrt{2 \times 10^2(1 - 0.50)} = 10.0000.$$

The formulas for matched designs apply to pre-post designs as well. The pre and post means correspond to the means in the matched groups,  $n$  is the number of subjects, and  $r$  is the correlation between pre-scores and post-scores.

### Calculation of effect size estimates from information that is reported

When a researcher has access to a full set of summary data such as the mean, standard deviation, and sample size for each group, the computation of the effect size and its variance is relatively straightforward. In practice, however, the researcher will often be working with only partial data. For example, a paper may publish only the  $p$ -value, means, and sample sizes from a test of significance, leaving it to the meta-analyst to back-compute the effect size and variance. For information on computing effect sizes from partial information, see Borenstein *et al.* (in preparation).

### Including different study designs in the same analysis

Sometimes a systematic review will include studies that used independent groups and also studies that used matched groups. From a statistical perspective the effect size ( $D$ ) has the same meaning regardless of the study design. Therefore, we can compute

the effect size and variance from each study using the appropriate formula, and then include all studies in the same analysis. While there is no technical barrier to using different study designs in the same analysis, there may be a concern that studies which used different designs might differ in substantive ways as well (see Chapter 45).

For all study designs (whether using independent or paired groups) the direction of the effect ( $\bar{X}_1 - \bar{X}_2$  or  $\bar{X}_2 - \bar{X}_1$ ) is arbitrary, except that the researcher must decide on a convention and then apply this consistently. For example, if a positive difference will indicate that the treated group did better than the control group, then this convention must apply for studies that used independent designs and for studies that used pre-post designs. In some cases it might be necessary to reverse the computed sign of the effect size to ensure that the convention is followed.

### STANDARDIZED MEAN DIFFERENCE, *d* AND *g*

As noted, the raw mean difference is a useful index when the measure is meaningful, either inherently or because of widespread use. By contrast, when the measure is less well known (for example, a proprietary scale with limited distribution), the use of a raw mean difference has less to recommend it. In any event, the raw mean difference is an option only if all the studies in the meta-analysis use the same scale. If different studies use different instruments (such as different psychological or educational tests) to assess the outcome, then the scale of measurement will differ from study to study and it would not be meaningful to combine raw mean differences.

In such cases we can divide the mean difference in each study by that study's standard deviation to create an index (the standardized mean difference) that would be comparable across studies. This is the same approach suggested by Cohen (1969, 1987) in connection with describing the magnitude of effects in statistical power analysis.

The standardized mean difference can be considered as being comparable across studies based on either of two arguments (Hedges & Olkin, 1985). If the outcome measures in all studies are linear transformations of each other, the standardized mean difference can be seen as the mean difference that would have been obtained if all data were transformed to a scale where the standard deviation within-groups was equal to 1.0.

The other argument for comparability of standardized mean differences is the fact that the standardized mean difference is a measure of overlap between distributions. In this telling, the standardized mean difference reflects the difference between the distributions in the two groups (and how each represents a distinct cluster of scores) even if they do not measure exactly the same outcome (see Cohen, 1987, Grissom & Kim, 2005).

Consider a study that uses two independent groups, and suppose we wish to compare the means of these two groups. Let  $\mu_1$  and  $\sigma_1$  be the true (population) mean and standard deviation of the first group and let  $\mu_2$  and  $\sigma_2$  be the true (population) mean and standard deviation of the other group. If the two population standard deviations

are the same (as is assumed in most parametric data analysis techniques), so that  $\sigma_1 = \sigma_2 = \sigma$ , then the standardized mean difference parameter or population standardized mean difference is defined as

$$\delta = \frac{\mu_1 - \mu_2}{\sigma}. \quad (4.17)$$

In the sections that follow, we show how to estimate  $\delta$  from studies that used independent groups, and from studies that used pre-post or matched group designs. It is also possible to estimate  $\delta$  from studies that used other designs (including clustered designs) but these are not addressed here (see resources at the end of this Part). We make the common assumption that  $\sigma_1^2 = \sigma_2^2$ , which allows us to pool the estimates of the standard deviation, and do not address the case where these are assumed to differ from each other.

### Computing $d$ and $g$ from studies that use independent groups

We can estimate the standardized mean difference ( $\delta$ ) from studies that used two independent groups as

$$d = \frac{\bar{X}_1 - \bar{X}_2}{S_{within}}. \quad (4.18)$$

In the numerator,  $\bar{X}_1$  and  $\bar{X}_2$  are the sample means in the two groups. In the denominator  $S_{within}$  is the within-groups standard deviation, pooled across groups,

$$S_{within} = \sqrt{\frac{(n_1 - 1)S_1^2 + (n_2 - 1)S_2^2}{n_1 + n_2 - 2}} \quad (4.19)$$

where  $n_1$  and  $n_2$  are the sample sizes in the two groups, and  $S_1$  and  $S_2$  are the standard deviations in the two groups. The reason that we pool the two sample estimates of the standard deviation is that even if we assume that the underlying population standard deviations are the same (that is  $\sigma_1 = \sigma_2 = \sigma$ ), it is unlikely that the sample estimates  $S_1$  and  $S_2$  will be identical. By pooling the two estimates of the standard deviation, we obtain a more accurate estimate of their common value.

The *sample estimate* of the standardized mean difference is often called Cohen's  $d$  in research synthesis. Some confusion about the terminology has resulted from the fact that the index  $\delta$ , originally proposed by Cohen as a *population parameter* for describing the size of effects for statistical power analysis, is also sometimes called  $d$ . In this volume we use the symbol  $\delta$  to denote the effect size parameter and  $d$  for the sample estimate of that parameter.

The variance of  $d$  is given (to a very good approximation) by

$$V_d = \frac{n_1 + n_2}{n_1 n_2} + \frac{d^2}{2(n_1 + n_2)}. \quad (4.20)$$

In this equation the first term on the right of the equals sign reflects uncertainty in the estimate of the mean difference (the numerator in (4.18)), and the second reflects uncertainty in the estimate of  $S_{within}$  (the denominator in (4.18)).

The standard error of  $d$  is the square root of  $V_d$ ,

$$SE_d = \sqrt{V_d}. \quad (4.21)$$

It turns out that  $d$  has a slight bias, tending to overestimate the absolute value of  $\delta$  in small samples. This bias can be removed by a simple correction that yields an unbiased estimate of  $\delta$ , with the unbiased estimate sometimes called Hedges'  $g$  (Hedges, 1981). To convert from  $d$  to Hedges'  $g$  we use a correction factor, which is called  $J$ . Hedges (1981) gives the exact formula for  $J$ , but in common practice researchers use an approximation,

$$J = 1 - \frac{3}{4df - 1}. \quad (4.22)$$

In this expression,  $df$  is the degrees of freedom used to estimate  $S_{within}$ , which for two independent groups is  $n_1 + n_2 - 2$ . This approximation always has error of less than 0.007 and less than 0.035 percent when  $df \geq 10$  (Hedges, 1981). Then,

$$g = J \times d \quad (4.23)$$

$$V_g = J^2 \times V_d, \quad (4.24)$$

and

$$SE_g = \sqrt{V_g}. \quad (4.25)$$

For example, suppose a study has sample means  $\bar{X}_1 = 103$ ,  $\bar{X}_2 = 100$ , sample standard deviations  $S_1 = 5.5$ ,  $S_2 = 4.5$ , and sample sizes  $n_1 = n_2 = 50$ . We would estimate the pooled-within-groups standard deviation as

$$S_{within} = \sqrt{\frac{(50 - 1) \times 5.5^2 + (50 - 1) \times 4.5^2}{50 + 50 - 2}} = 5.0249.$$

Then,

$$d = \frac{103 - 100}{5.0249} = 0.5970,$$

$$V_d = \frac{50 + 50}{50 \times 50} + \frac{0.5970^2}{2(50 + 50)} = 0.0418,$$

and

$$SE_d = \sqrt{0.0418} = 0.2044.$$

The correction factor ( $J$ ), Hedges'  $g$ , its variance and standard error are given by

$$J = \left(1 - \frac{3}{4 \times 98 - 1}\right) = 0.9923,$$

$$g = 0.9923 \times 0.5970 = 0.5924,$$

$$V_g = 0.9923^2 \times 0.0418 = 0.0411,$$

and

$$SE_g = \sqrt{0.0411} = 0.2028.$$

The correction factor ( $J$ ) is always less than 1.0, and so  $g$  will always be less than  $d$  in absolute value, and the variance of  $g$  will always be less than the variance of  $d$ . However,  $J$  will be very close to 1.0 unless  $df$  is very small (say, less than 10) and so (as in this example) the difference is usually trivial (Hedges, 1981).

Some slightly different expressions for the variance of  $d$  (and  $g$ ) have been given by different authors and even the same authors at different times. For example, the denominator of the second term of the variance of  $d$  is given here as  $2(n_1 + n_2)$ . This expression is obtained by one method (assuming the  $n$ 's become large with  $\delta$  fixed). An alternate derivation (assuming  $n$ 's become large with  $\sqrt{n}\delta$  fixed) leads to a denominator in the second term that is slightly different, namely  $2(n_1 + n_2 - 2)$ . Unless  $n_1$  and  $n_2$  are very small, these expressions will be almost identical.

Similarly, the expression given here for the variance of  $g$  is  $J^2$  times the variance of  $d$ , but many authors ignore the  $J^2$  term because it is so close to unity in most cases. Again, while it is preferable to include this correction factor, the inclusion of this factor is likely to make little practical difference.

### Computing $d$ and $g$ from studies that use pre-post scores or matched groups

We can estimate the standardized mean difference ( $\delta$ ) from studies that used matched groups or pre-post scores in one group. The formula for the sample estimate of  $d$  is

$$d = \frac{\bar{Y}_{\text{diff}}}{S_{\text{within}}} = \frac{\bar{Y}_1 - \bar{Y}_2}{S_{\text{within}}}. \quad (4.26)$$

This is the same formula as for independent groups (4.18). However, when we are working with independent groups, the natural unit of deviation is the standard deviation within groups and so this value is typically reported (or easily imputed). By contrast, when we are working with matched groups, the natural unit of deviation is the standard deviation of the difference scores, and so *this* is the value that is likely to be reported. To compute  $d$  from the standard deviation of the differences we need to impute the standard deviation within groups, which would then serve as the denominator in (4.26).

Concretely, when working with a matched study, the standard deviation within groups can be imputed from the standard deviation of the difference, using

$$S_{\text{within}} = \frac{S_{\text{diff}}}{\sqrt{2(1 - r)}}, \quad (4.27)$$

where  $r$  is the correlation between pairs of observations (e.g., the pretest-posttest correlation). Then we can apply (4.26) to compute  $d$ . The variance of  $d$  is given by

$$V_d = \left( \frac{1}{n} + \frac{d^2}{2n} \right) 2(1 - r), \quad (4.28)$$

where  $n$  is the number of pairs. The standard error of  $d$  is just the square root of  $V_d$ ,

$$SE_d = \sqrt{V_d}. \quad (4.29)$$

Since the correlation between pre- and post-scores is required to impute the standard deviation within groups from the standard deviation of the difference, we must assume

that this correlation is known or can be estimated with high precision. Otherwise we may estimate the correlation from related studies, and possibly perform a sensitivity analysis using a range of plausible correlations.

To compute Hedges'  $g$  and associated statistics we would use formulas (4.22) through (4.25). The degrees of freedom for computing  $J$  is  $n - 1$ , where  $n$  is the number of pairs.

For example, suppose that a study has pre-test and post-test sample means  $\bar{X}_1 = 103$ ,  $\bar{X}_2 = 100$ , sample standard deviation of the difference  $S_{diff} = 5.5$ , sample size  $n = 50$ , and a correlation between pre-test and post-test of  $r = 0.7$ . The standard deviation within groups is imputed from the standard deviation of the difference by

$$S_{within} = \frac{5.5}{\sqrt{2(1 - 0.7)}} = 7.1005.$$

Then  $d$ , its variance and standard error are computed as

$$d = \frac{103 - 100}{7.1000} = 0.4225,$$

$$v_d = \left( \frac{1}{50} + \frac{0.4225^2}{2 \times 50} \right) (2(1 - 0.7)) = 0.0131,$$

and

$$SE_d = \sqrt{0.0131} = 0.1143.$$

The correction factor  $J$ , Hedges'  $g$ , its variance and standard error are given by

$$J = \left( 1 - \frac{3}{4 \times 49 - 1} \right) = 0.9846,$$

$$g = 0.9846 \times 0.4225 = 0.4160,$$

$$V_g = 0.9846^2 \times 0.0131 = 0.0127,$$

and

$$SE_g = \sqrt{0.0127} = 0.1126.$$

### Including different study designs in the same analysis

As we noted earlier, a single systematic review can include studies that used independent groups and also studies that used matched groups. From a statistical perspective the effect size ( $d$  or  $g$ ) has the same meaning regardless of the study design. Therefore, we can compute the effect size and variance from each study using the appropriate formula, and then include all studies in the same analysis. While there are no technical barriers to using studies with different designs in the same analysis, there may be a concern that these studies could differ in substantive ways as well (see Chapter 45).

For all study designs the direction of the effect ( $\bar{X}_1 - \bar{X}_2$  or  $\bar{X}_2 - \bar{X}_1$ ) is arbitrary, except that the researcher must decide on a convention and then apply this consistently. For example, if a positive difference indicates that the treated group did better than the control group, then this convention must apply for studies that used independent

designs and for studies that used pre-post designs. It must also apply for all outcome measures. In some cases (for example, if some studies defined outcome as the number of correct answers while others defined outcome as the number of mistakes) it will be necessary to reverse the computed sign of the effect size to ensure that the convention is applied consistently.

## RESPONSE RATIOS

In research domains where the outcome is measured on a physical scale (such as length, area, or mass) and is unlikely to be zero, the ratio of the means in the two groups might serve as the effect size index. In experimental ecology this effect size index is called the response ratio (Hedges, Gurevitch, & Curtis, 1999). It is important to recognize that the response ratio is only meaningful when the outcome is measured on a true ratio scale. The response ratio is not meaningful for studies (such as most social science studies) that measure outcomes such as test scores, attitude measures, or judgments, since these have no natural scale units and no natural zero points.

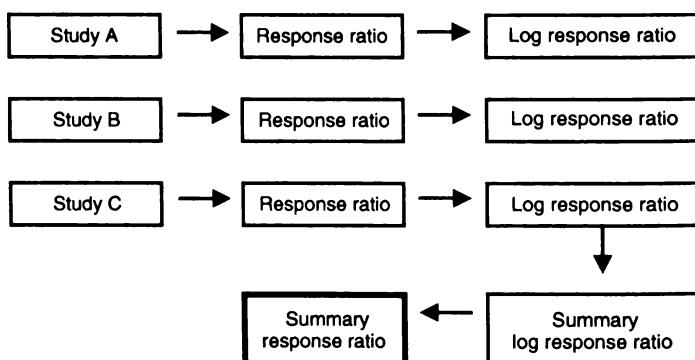
For response ratios, computations are carried out on a log scale (see the discussion under risk ratios, below, for an explanation). We compute the log response ratio and the standard error of the log response ratio, and use these numbers to perform all steps in the meta-analysis. Only then do we convert the results back into the original metric. This is shown schematically in Figure 4.1.

The response ratio is computed as

$$R = \frac{\bar{X}_1}{\bar{X}_2} \quad (4.30)$$

where  $\bar{X}_1$  is the mean of group 1 and  $\bar{X}_2$  is the mean of group 2. The log response ratio is computed as

$$\ln R = \ln(R) = \ln\left(\frac{\bar{X}_1}{\bar{X}_2}\right) = \ln(\bar{X}_1) - \ln(\bar{X}_2). \quad (4.31)$$



**Figure 4.1** Response ratios are analyzed in log units.

The variance of the log response ratio is approximately

$$V_{\ln R} = S_{pooled}^2 \left( \frac{1}{n_1(\bar{X}_1)^2} + \frac{1}{n_2(\bar{X}_2)^2} \right), \quad (4.32)$$

where  $S_{pooled}$  is the pooled standard deviation. The approximate standard error is

$$SE_{\ln R} = \sqrt{V_{\ln R}}. \quad (4.33)$$

Note that we do not compute a variance for the response ratio in its original metric. Rather, we use the *log* response ratio and its variance in the analysis to yield a summary effect, confidence limits, and so on, in log units. We then convert each of these values back to response ratios using

$$R = \exp(\ln R), \quad (4.34)$$

$$LL_R = \exp(LL_{\ln R}), \quad (4.35)$$

and

$$UL_R = \exp(UL_{\ln R}), \quad (4.36)$$

where  $LL$  and  $UL$  represent the lower and upper limits, respectively.

For example, suppose that a study has two independent groups with means  $\bar{X}_1 = 61.515$ ,  $\bar{X}_2 = 51.015$ , pooled within-group standard deviation 19.475, and sample size  $n_1 = n_2 = 10$ .

Then  $R$ , its variance and standard error are computed as

$$R = \frac{61.515}{51.015} = 1.2058,$$

$$\ln R = \ln(1.2058) = 0.1871,$$

$$V_{\ln R} = 19.475^2 \left( \frac{1}{10 \times (61.515)^2} + \frac{1}{10 \times (51.015)^2} \right) = 0.0246.$$

and

$$SE_{\ln R} = \sqrt{0.0246} = 0.1581.$$

### SUMMARY POINTS

- The raw mean difference ( $D$ ) may be used as the effect size when the outcome scale is either inherently meaningful or well known due to widespread use. This effect size can only be used when all studies in the analysis used precisely the same scale.
- The standardized mean difference ( $d$  or  $g$ ) transforms all effect sizes to a common metric, and thus enables us to include different outcome measures in the same synthesis. This effect size is often used in primary research as well as meta-analysis, and therefore will be intuitive to many researchers.

- The response ratio ( $R$ ) is often used in ecology. This effect size is only meaningful when the outcome has a natural zero point, but when this condition holds, it provides a unique perspective on the effect size.
- It is possible to compute an effect size and variance from studies that used two independent groups, from studies that used matched groups (or pre-post designs) and from studies that used clustered groups. These effect sizes may then be included in the same meta-analysis.

# Effect Sizes Based on Binary Data ( $2 \times 2$ Tables)

---

- Introduction
  - Risk ratio
  - Odds ratio
  - Risk difference
  - Choosing an effect size index
- 

## INTRODUCTION

For data from a prospective study, such as a randomized trial, that was originally reported as the number of events and non-events in two groups (the classic  $2 \times 2$  table), researchers typically compute a risk ratio, an odds ratio, and/or a risk difference. This data can be represented as cells A, B, C, and D, as shown in Table 5.1.

For example, assume a study with a sample size of 100 per group. Five patients died in the treated group, as compared with ten who died in the control group (see Table 5.2).

From these data we might compute a risk ratio, an odds ratio, and/or a risk difference.

## RISK RATIO

The risk ratio is simply the ratio of two *risks*. Here, the risk of death in the treated group is 5/100 and the risk of death in the control group is 10/100, so the ratio of the two risks is 0.50. This index has the advantage of being intuitive, in the sense that the meaning of a ratio is clear.

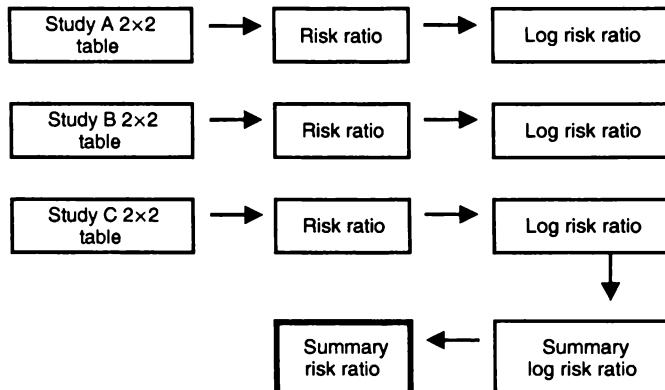
For risk ratios, computations are carried out on a log scale. We compute the log risk ratio, and the standard error of the log risk ratio, and will use these numbers to perform all steps in the meta-analysis. Only then will we convert the results back into the original metric. This is shown schematically in Figure 5.1.

**Table 5.1** Nomenclature for  $2 \times 2$  table of outcome by treatment.

	Events	Non-Events	N
Treated	A	B	$n_1$
Control	C	D	$n_2$

**Table 5.2** Fictional data for a  $2 \times 2$  table.

	Dead	Alive	N
Treated	5	95	100
Control	10	90	100

**Figure 5.1** Risk ratios are analyzed in log units.

The computational formula for the risk ratio is

$$\text{RiskRatio} = \frac{A/n_1}{C/n_2}. \quad (5.1)$$

The log risk ratio is then

$$\text{LogRiskRatio} = \ln(\text{RiskRatio}), \quad (5.2)$$

with approximate variance

$$V_{\text{LogRiskRatio}} = \frac{1}{A} - \frac{1}{n_1} + \frac{1}{C} - \frac{1}{n_2}, \quad (5.3)$$

and approximate standard error

$$SE_{\text{LogRiskRatio}} = \sqrt{V_{\text{LogRiskRatio}}}. \quad (5.4)$$

Note that we do not compute a variance for the risk ratio in its original metric. Rather, we use the *log* risk ratio and its variance in the analysis to yield a summary effect,

confidence limits, and so on, in log units. We then convert each of these values back to risk ratios using

$$\text{RiskRatio} = \exp(\text{LogRiskRatio}), \quad (5.5)$$

$$\text{LL}_{\text{RiskRatio}} = \exp(\text{LL}_{\text{LogRiskRatio}}), \quad (5.6)$$

and

$$\text{UL}_{\text{RiskRatio}} = \exp(\text{UL}_{\text{LogRiskRatio}}), \quad (5.7)$$

where  $\text{LL}$  and  $\text{UL}$  represent the lower and upper limits, respectively.

In the running example the risk ratio is

$$\text{RiskRatio} = \frac{5/100}{10/100} = 0.5000.$$

The log is

$$\text{LogRiskRatio} = \ln(0.5000) = -0.6932,$$

with variance

$$V_{\text{LogRiskRatio}} = \frac{1}{5} - \frac{1}{100} + \frac{1}{10} - \frac{1}{100} = 0.2800,$$

and standard error

$$SE_{\text{LogRiskRatio}} = \sqrt{0.280} = 0.5292.$$

Note 1. The log transformation is needed to maintain symmetry in the analysis. Assume that one study reports that the risk is twice as high in Group A while another reports that it is twice as high in Group B. Assuming equal weights, these studies should balance each other, with a combined effect showing equal risks (a risk ratio of 1.0). However, on the ratio scale these correspond to risk ratios of 0.50 and 2.00, which would yield a mean of 1.25. By working with log values we can avoid this problem. In log units the two estimates are  $-0.693$  and  $+0.693$ , which yield a mean of 0.00. We convert this back to a risk ratio of 1.00, which is the correct value for this data.

Note 2. Although we defined the risk ratio in this example as

$$\text{RiskRatio} = \frac{5/100}{10/100} = 0.5000$$

(which gives the risk ratio of dying) we could alternatively have focused on the *risk of staying alive*, given by

$$\text{RiskRatio} = \frac{95/100}{90/100} = 1.0556.$$

The ‘risk’ of staying alive is *not* the inverse of the risk of dying (that is, 1.056 is not the inverse of 0.50), and therefore this should be considered a different measure of effect size.

## ODDS RATIO

Where the risk ratio is the ratio of two *risks*, the odds ratio is the ratio of two *odds*. Here, the odds of death in the treated group would be 5/95, or 0.0526 (since probability of death in the treated group is 5/100 and the probability of life is 95/100), while the

odds of death in the control group would be 10/90, or 0.1111. The ratio of the two odds would then be 0.0526/0.1111, or 0.4737.

Many people find this effect size measure less intuitive than the risk ratio, but the odds ratio has statistical properties that often make it the best choice for a meta-analysis. When the risk of the event is low, the odds ratio will be similar to the risk ratio.

For odds ratios, computations are carried out on a log scale (for the same reason as for risk ratios). We compute the log odds ratio, and the standard error of the log odds ratio, and will use these numbers to perform all steps in the meta-analysis. Only then will we convert the results back into the original metric. This is shown schematically in Figure 5.2.

The computational formula for the odds ratio is

$$\text{OddsRatio} = \frac{AD}{BC}. \quad (5.8)$$

The log odds ratio is then

$$\text{LogOddsRatio} = \ln(\text{OddsRatio}), \quad (5.9)$$

with approximate variance

$$V_{\text{LogOddsRatio}} = \frac{1}{A} + \frac{1}{B} + \frac{1}{C} + \frac{1}{D} \quad (5.10)$$

and approximate standard error

$$SE_{\text{LogOddsRatio}} = \sqrt{V_{\text{LogOddsRatio}}}. \quad (5.11)$$

Note that we do not compute a variance for the odds ratio. Rather, the log odds ratio and its variance are used in the analysis to yield a summary effect, confidence limits, and so on, in log units. We then convert each of these values back to odds ratios using

$$\text{OddsRatio} = \exp(\text{LogOddsRatio}), \quad (5.12)$$

$$LL_{\text{OddsRatio}} = \exp(LL_{\text{LogOddsRatio}}), \quad (5.13)$$

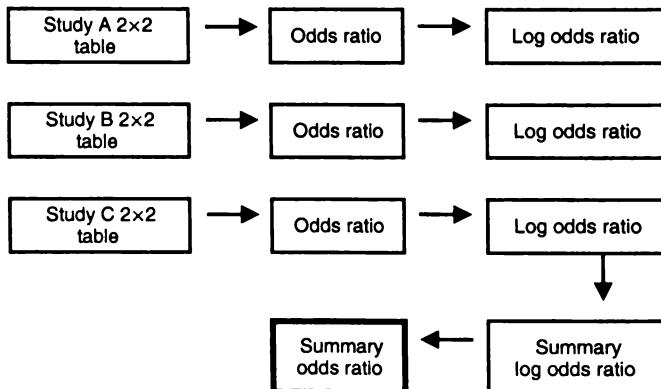


Figure 5.2 Odds ratios are analyzed in log units.

and

$$UL_{OddsRatio} = \exp(UL_{LogOddsRatio}), \quad (5.14)$$

where  $LL$  and  $UL$  represent the lower and upper limits, respectively.

In the running example

$$OddsRatio = \frac{5 \times 90}{95 \times 10} = 0.4737$$

and

$$LogOddsRatio = \ln(0.4737) = -0.7472,$$

with variance

$$V_{LogOddsRatio} = \frac{1}{5} + \frac{1}{95} + \frac{1}{10} + \frac{1}{90} = 0.3216$$

and standard error

$$SE_{LogOddsRatio} = \sqrt{0.3216} = 0.5671.$$

Note. When working with the odds ratio or risk ratio we can place either the Treated group or the Control group in the numerator, as long we apply this consistently across all studies. If we put the Treated group in the denominator the log odds ratio would change signs (from  $-0.7472$  to  $+0.7472$ ) and the odds ratio would change to its inverse (from  $0.4737$  to  $2.1110$ ). The same thing happens to the odds ratio if we swap Dead and Alive within each group. However, this is *not* the case for the risk ratio.

## RISK DIFFERENCE

The risk difference is the *difference* between two risks. Here, the risk in the treated group is 0.05 and the risk in the control group is 0.10, so the risk difference is  $-0.05$ .

Unlike the case for risk ratios and for odds ratios, computations for risk differences are carried out in raw units rather than log units.

The risk difference is defined as

$$RiskDiff = \left( \frac{A}{n_1} \right) - \left( \frac{C}{n_2} \right) \quad (5.15)$$

with approximate variance

$$V_{RiskDiff} = \frac{AB}{n_1^3} + \frac{CD}{n_2^3} \quad (5.16)$$

and approximate standard error

$$SE_{RiskDiff} = \sqrt{V_{RiskDiff}}. \quad (5.17)$$

In the running example

$$RiskDiff = \left( \frac{5}{100} \right) - \left( \frac{10}{100} \right) = -0.0500$$

with variance

$$V_{RiskDiff} = \frac{5 \times 95}{100^3} + \frac{10 \times 90}{100^3} = 0.0014$$

and standard error

$$SE_{RiskDiff} = \sqrt{0.00138} = 0.0371.$$

## CHOOSING AN EFFECT SIZE INDEX

In selecting among the risk ratio, odds ratio, and risk difference the researcher needs to consider both substantive and technical factors.

The risk ratio and odds ratio are relative measures, and therefore tend to be relatively insensitive to differences in baseline events. By contrast, the risk difference is an absolute measure and as such is very sensitive to the baseline risk. If we wanted to test a compound and believed that it reduced the risk of an event by 20% regardless of the baseline risk, then by using a ratio index we would expect to see the same effect size across studies even if the baseline risk varied from study to study. The risk difference, by contrast, would be higher in studies with a higher base rate.

At the same time, if we wanted to convey the clinical impact of the treatment, the risk difference might be the better measure. Suppose we perform a meta-analysis to assess the risk of adverse events for treated versus control groups. The risk is 1/1000 for treated patients versus 1/2000 for control patients, for a risk ratio of 2.00. At the same time, the risk difference is 0.0010 versus 0.0005 for a risk difference of 0.0005. These two numbers (2.00 and 0.0005) are both correct, but measure different things.

Because the ratios are less sensitive to baseline risk while the risk difference is sometimes more clinically meaningful, some suggest using the risk ratio (or odds ratio) to perform the meta-analysis and compute a summary risk (or odds) ratio. Then, they can use this to predict the risk difference for any given baseline risk.

### SUMMARY POINTS

- We can compute the risk of an event (such as the risk of death) in each group (for example, treated versus control). The ratio of these risks then serves as an effect size (the risk ratio).
- We can compute the odds of an event (such as ratio of dying to living) in each group (for example, treated versus control). The ratio of these odds then serves as the odds ratio.
- We can compute the risk of an event (such as the risk of death) in each group (for example, treated versus control). The difference in these risks then serves as an effect size (the risk difference).
- To work with the risk ratio or odds ratio we transform all values to log values, perform the analyses, and then convert the results back to ratio values for presentation. To work with the risk difference we work with the raw values.

# Effect Sizes Based on Correlations

---

- Introduction
  - Computing  $r$
  - Other approaches
- 

## INTRODUCTION

For studies that report a correlation between two continuous variables, the correlation coefficient itself can serve as the effect size index. The correlation is an intuitive measure that, like  $\delta$ , has been standardized to take account of different metrics in the original scales. The population parameter is denoted by  $\rho$  (the Greek letter rho).

## COMPUTING $r$

The estimate of the correlation parameter  $\rho$  is simply the sample correlation coefficient,  $r$ . The variance of  $r$  is approximately

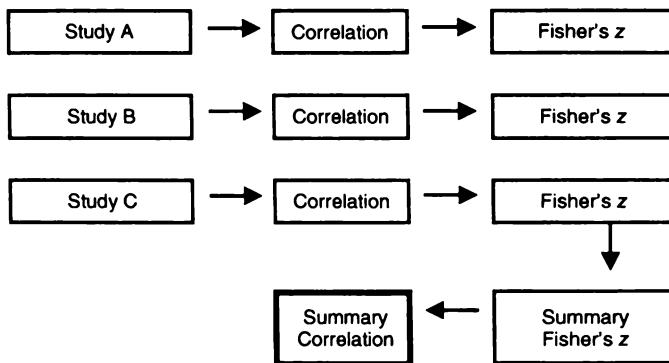
$$V_r = \frac{(1 - r^2)^2}{n - 1} \quad (6.1)$$

where  $n$  is the sample size.

Most meta-analysts do not perform syntheses on the correlation coefficient itself because the variance depends strongly on the correlation. Rather, the correlation is converted to the Fisher's  $z$  scale (not to be confused with the  $z$ -score used with significance tests), and all analyses are performed using the transformed values. The results, such as the summary effect and its confidence interval, would then be converted back to correlations for presentation. This is shown schematically in Figure 6.1, and is analogous to the procedure used with odds ratios or risk ratios where all analyses are performed using log transformed values, and then converted back to the original metric.

The transformation from sample correlation  $r$  to Fisher's  $z$  is given by

$$z = 0.5 \times \ln \left( \frac{1 + r}{1 - r} \right). \quad (6.2)$$



**Figure 6.1** Correlations are analyzed in Fisher's  $z$  units.

The variance of  $z$  (to an excellent approximation) is

$$V_z = \frac{1}{n - 3}, \quad (6.3)$$

and the standard error is

$$SE_z = \sqrt{V_z}. \quad (6.4)$$

When working with Fisher's  $z$ , we do not use the variance for the correlation. Rather, the Fisher's  $z$  score and its variance are used in the analysis, which yield a summary effect, confidence limits, and so on, in the Fisher's  $z$  metric. We then convert each of these values back to correlation units using

$$r = \frac{e^{2z} - 1}{e^{2z} + 1}. \quad (6.5)$$

For example, if a study reports a correlation of 0.50 with a sample size of 100, we would compute

$$z = 0.5 \times \ln \left( \frac{1 + 0.5}{1 - 0.5} \right) = 0.5493,$$

$$V_z = \frac{1}{100 - 3} = 0.0103,$$

and

$$SE_z = \sqrt{0.0103} = 0.1015.$$

To convert the Fisher's  $z$  value back to a correlation, we would use

$$r = \frac{e^{(2 \times 0.5493)} - 1}{e^{(2 \times 0.5493)} + 1} = 0.5000.$$

## OTHER APPROACHES

Hunter and Schmidt (2004) advocate an approach for working with correlations that differs in several ways from the one presented here. This approach is discussed in Chapter 38.

**SUMMARY POINTS**

- When studies report data as correlations, we usually use the correlation coefficient itself as the effect size. We transform the correlation using the Fisher's z transformation and perform the analysis using this index. Then, we convert the summary values back to correlations for presentation.



# Converting Among Effect Sizes

---

### Introduction

Converting from the log odds ratio to  $d$

Converting from  $d$  to the log odds ratio

Converting from  $r$  to  $d$

Converting from  $d$  to  $r$

---

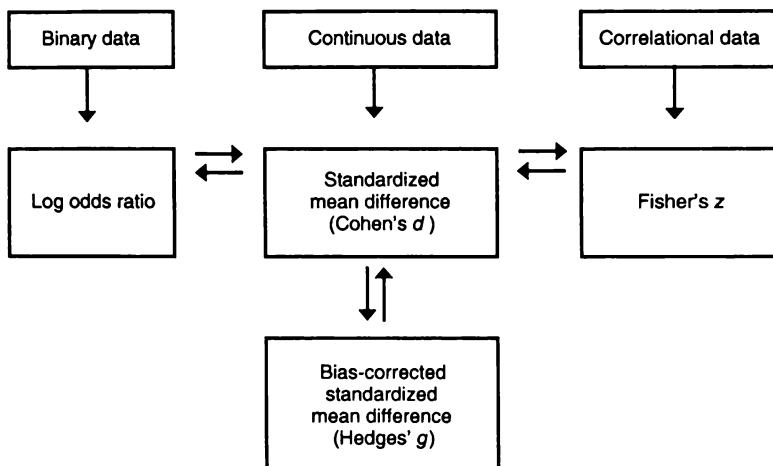
## INTRODUCTION

Earlier in this Part we discussed the case where different study designs were used to compute the same effect size. For example, studies that used independent groups and studies that used matched groups were both used to yield estimates of the standardized mean difference,  $g$ . There is no problem in combining these estimates in a meta-analysis since the effect size has the same meaning in all studies.

Consider, however, the case where some studies report a difference in means, which is used to compute a standardized mean difference; others report a difference in proportions, which is used to compute an odds ratio; and others report a correlation. All the studies address the same broad question, and we want to include them in one meta-analysis. Unlike the earlier case, we are now dealing with different indices, and we need to convert them to a common index before we can proceed.

The question of whether or not it is appropriate to combine effect sizes from studies that used different metrics must be considered on a case by case basis. The key issue is that it only makes sense to compute a summary effect from studies that we judge to be comparable in relevant ways. If we would be comfortable combining these studies if they had used the same metric, then the fact that they used different metrics should not be an impediment.

For example, suppose that several randomized controlled trials start with the same measure, on a continuous scale, but some report the outcome as a mean and others dichotomize the outcome and report it as success or failure. In this case, it may be highly appropriate to transform the standardized mean differences and the odds ratios to a common metric and then combine them across studies. By contrast, observational



**Figure 7.1** Converting among effect sizes.

studies that report correlations may be substantially different from observational studies that report odds ratios. In this case, even if there is no technical barrier to converting the effects to a common metric, it may be a bad idea from a substantive perspective.

In this chapter we present formulas for converting between an odds ratio and  $d$ , or between  $d$  and  $r$ . By combining formulas it is also possible to convert from an odds ratio, via  $d$ , to  $r$  (see Figure 7.1). In every case the formula for converting the effect size is accompanied by a formula to convert the variance.

When we convert between different measures we make certain assumptions about the nature of the underlying traits or effects. Even if these assumptions do not hold exactly, the decision to use these conversions is often better than the alternative, which is to simply omit the studies that happened to use an alternate metric. This would involve loss of information, and possibly the *systematic* loss of information, resulting in a biased sample of studies. A sensitivity analysis to compare the meta-analysis results with and without the converted studies would be important.

Figure 7.1 outlines the mechanism for incorporating multiple kinds of data in the same meta-analysis. First, each study is used to compute an effect size and variance of its native index, the log odds ratio for binary data,  $d$  for continuous data, and  $r$  for correlational data. Then, we convert all of these indices to a common index, which would be either the log odds ratio,  $d$ , or  $r$ . If the final index is  $d$ , we can move from there to Hedges'  $g$ . This common index and its variance are then used in the analysis.

## CONVERTING FROM THE LOG ODDS RATIO TO $d$

We can convert from a log odds ratio (*LogOddsRatio*) to the standardized mean difference  $d$  using

$$d = \text{LogOddsRatio} \times \frac{\sqrt{3}}{\pi}, \quad (7.1)$$

where  $\pi$  is the mathematical constant (approximately 3.14159). The variance of  $d$  would then be

$$V_d = V_{\text{LogOddsRatio}} \times \frac{3}{\pi^2}, \quad (7.2)$$

where  $V_{\text{LogOddsRatio}}$  is the variance of the log odds ratio. This method was originally proposed by Hasselblad and Hedges (1995) but variations have been proposed (see Sanchez-Meca, Marin-Martinez, & Chacon-Moscoso, 2003; Whitehead, 2002). It assumes that an underlying continuous trait exists and has a logistic distribution (which is similar to a normal distribution) in each group. In practice, it will be difficult to test this assumption.

For example, if the log odds ratio were  $\text{LogOddsRatio} = 0.9069$  with a variance of  $V_{\text{LogOddsRatio}} = 0.0676$ , then

$$d = 0.9069 \times \frac{\sqrt{3}}{3.1416} = 0.5000$$

with variance

$$V_d = 0.0676 \times \frac{3}{3.1416^2} = 0.0205.$$

## CONVERTING FROM $d$ TO THE LOG ODDS RATIO

We can convert from the standardized mean difference  $d$  to the log odds ratio ( $\text{LogOddsRatio}$ ) using

$$\text{LogOddsRatio} = d \frac{\pi}{\sqrt{3}}, \quad (7.3)$$

where  $\pi$  is the mathematical constant (approximately 3.14159). The variance of  $\text{LogOddsRatio}$  would then be

$$V_{\text{LogOddsRatio}} = V_d \frac{\pi^2}{3}. \quad (7.4)$$

For example, if  $d = 0.5000$  and  $V_d = 0.0205$  then

$$\text{log OddsRatio} = 0.5000 \times \frac{3.1416}{\sqrt{3}} = 0.9069,$$

and

$$V_{\text{LogOddsRatio}} = 0.0205 \times \frac{3.1416^2}{3} = 0.0676.$$

To employ this transformation we assume that the continuous data have the logistic distribution.

## CONVERTING FROM $r$ TO $d$

We convert from a correlation ( $r$ ) to a standardized mean difference ( $d$ ) using

$$d = \frac{2r}{\sqrt{1 - r^2}}. \quad (7.5)$$

The variance of  $d$  computed in this way (converted from  $r$ ) is

$$V_d = \frac{4V_r}{(1 - r^2)^3} \quad (7.6)$$

For example, if  $r = 0.50$  and  $V_r = 0.0058$ , then

$$d = \frac{2 \times 0.50}{\sqrt{1 - 0.50^2}} = 1.1547$$

and the variance of  $d$  is

$$V_d = \frac{4 \times 0.0058}{(1 - 0.50^2)^3} = 0.0550.$$

In applying this conversion we assume that the continuous data used to compute  $r$  has a bivariate normal distribution and that the two groups are created by dichotomizing one of the two variables.

## CONVERTING FROM $d$ TO $r$

We can convert from a standardized mean difference ( $d$ ) to a correlation ( $r$ ) using

$$r = \frac{d}{\sqrt{d^2 + a}} \quad (7.7)$$

where  $a$  is a correction factor for cases where  $n_1 \neq n_2$

$$a = \frac{(n_1 + n_2)^2}{n_1 n_2}. \quad (7.8)$$

The correction factor ( $a$ ) depends on the ratio of  $n_1$  to  $n_2$ , rather than the absolute values of these numbers. Therefore, if  $n_1$  and  $n_2$  are not known precisely, use  $n_1 = n_2$ , which will yield  $a = 4$ . The variance of  $r$  computed in this way (converted from  $d$ ) is

$$V_r = \frac{a^2 V_d}{(d^2 + a)^3}. \quad (7.9)$$

For example, if  $n_1 = n_2$ ,  $d = 1.1547$  and  $V_d = 0.0550$ , then

$$r = \frac{1.1547}{\sqrt{1.1547^2 + 4}} = 0.5000$$

and the variance of  $r$  converted from  $d$  will be

$$V_r = \frac{4^2 \times 0.0550}{(1.1547^2 + 4)^3} = 0.0058.$$

In applying this conversion, assume that a continuous variable was dichotomized to create the treatment and control groups.

When we transform between Fisher's  $z$  and  $d$  we are making assumptions about the independent variable only. When we transform between the log odds ratio and  $d$  we are making assumptions about the dependent variable only. As such, the two sets of assumptions are independent of each other, and one has no implications for the validity of the other. Therefore, we can apply both sets of assumptions and transform from Fisher's  $z$  through  $d$  to the log odds ratio, as well as the reverse.

### SUMMARY POINTS

- If all studies in the analysis are based on the same kind of data (means, binary, or correlational), the researcher should select an effect size based on that kind of data.
- When some studies use means, others use binary data, and others use correlational data, we can apply formulas to convert among effect sizes.
- Studies that used different measures may differ from each other in substantive ways, and we need to consider this possibility when deciding if it makes sense to include the various studies in the same analysis.



# Factors that Affect Precision

---

Introduction

Factors that affect precision

Sample size

Study design

---

## INTRODUCTION

In the preceding chapters we showed how to compute the variance for specific effect sizes such as the standardized mean difference or a log risk ratio. Our goal in this chapter is to provide some context for those formulas.

We use the term precision as a general term to encompass three formal statistics, the variance, standard error, and confidence interval. These are all related to each other, so when we discuss the impact of a factor on precision, this translates into an impact on all three. In this chapter we outline the relationship between these three indices of precision. Then, we discuss two factors that affect precision and make some studies more precise than others.

### Variance, standard error, and confidence intervals

The word “variance” can refer either to the population variance or the error variance (the variance of the mean). In this chapter we intend the latter meaning.

The variance is a measure of the mean squared deviation from the mean effect. For an effect size  $Y$  (used generically), the variance would be denoted simply as

$$V_Y. \tag{8.1}$$

The computation of the variance is different for every effect size index (some formulas were presented in the preceding chapters).

The variance has properties that make it useful for some statistical computations, but because its metric is based on squared values it is not an intuitive index. A more accessible index is the standard error, which is on the same scale as the effect size

itself. If  $Y$  is the effect size and  $V_Y$  is the variance of  $Y$ , then the standard error of  $Y$  ( $SE_Y$ ) is given by

$$SE_Y = \sqrt{V_Y}. \quad (8.2)$$

If we assume that the effect size is normally distributed then we can compute a 95% confidence interval using

$$LL_Y = \bar{Y} - 1.96 \times SE_Y \quad (8.3)$$

and

$$UL_Y = \bar{Y} + 1.96 \times SE_Y. \quad (8.4)$$

In these equations 1.96 is the  $Z$ -value corresponding to confidence limits of 95% (allowing for 2.5% error at either end of the distribution). We can also compute a test statistic  $Z$  as

$$Z_Y = \frac{\bar{Y}}{SE_Y}. \quad (8.5)$$

There is a relationship between the  $p$ -value for  $Z$  and the confidence interval, such that (in almost all cases) the  $p$ -value will be less than 0.05 if and only if the confidence interval does not include the null value.

## FACTORS THAT AFFECT PRECISION

Some of the factors that affect precision are unique to each effect size index, as explained in the preceding chapters. They are also unique to each study since each study has inherent factors, such as the homogeneity of the sample, which affect precision. Beyond these unique factors, however, are two factors that have an important and predictable impact on precision. One is the size of the sample and the other is the study design (whether the study used paired groups, independent groups, or clustered groups). The impact of these two factors is explained here.

## SAMPLE SIZE

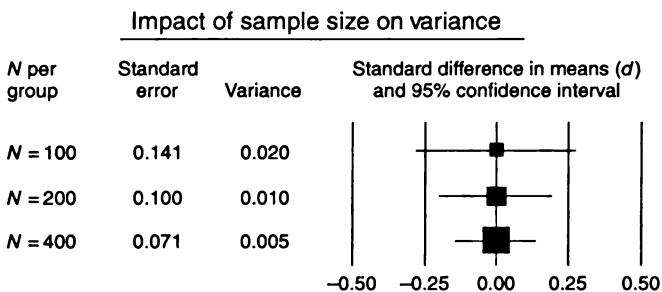
A dominant factor in precision is the sample size, with larger samples yielding more precise estimates than smaller samples.

For example, consider the three studies in Table 8.1. These studies compared the means in two independent groups, and we computed the standardized mean difference ( $d$ ), which is 0.0 in this example. The sample sizes in the three studies ( $A$ ,  $B$ ,  $C$ ) are 100, 200, and 400 per group, respectively, and the variances are 0.020, 0.010, and 0.005. In other words, as the sample size increases by a factor of 4 (compare studies  $A$  and  $C$ ) the variance will decrease by a factor of 4 and the standard error will decrease by a factor of 2 (that is, by the square root of 4).

Note. In this example we assume that  $d = 0.0$  which allows us to focus on the relationship between sample size and variance. When  $d$  is nonzero,  $d$  has an impact on the variance (though this impact is typically small).

**Table 8.1** Impact of sample size on variance.

Study	Design	N per group	Standard error	Variance
A	Independent	100	0.141	0.020
B	Independent	200	0.100	0.010
C	Independent	400	0.071	0.005

**Figure 8.1** Impact of sample size on variance.

The same information is presented graphically in Figure 8.1, where each study is represented by a box and bounded by a confidence interval. In this figure:

- The area of a box is proportional to the inverse of that study's variance.
- Any side on a box is proportional to the inverse of that study's standard error.
- The confidence interval for each box is proportional to that study's standard error.

Later, we will discuss how weights are assigned to each study in the meta-analysis. Under one scheme weights are inversely proportional to the variance, and study C would be assigned four times as much weight as study A.

## STUDY DESIGN

In the preceding example (where we compared different sample sizes) we assumed that the studies used two independent groups. Here, we consider what happens if we use a comparable sample size but an alternate study design.

One alternate design is matched pairs, where each person in the treated group is matched with a similar person in the control group (say, a sibling, or a person at the same disease stage). This design allows us to work with differences within these pairs (rather than differences between groups) which can reduce the error term and thus increase the precision of the estimate. The impact on precision depends on the correlation between (for example) siblings, with a higher correlation yielding greater precision.

In Table 8.2 line D (Independent) shows the variance for a study with 100 subjects per group, and is identical to Study A in the prior table. The three lines below this

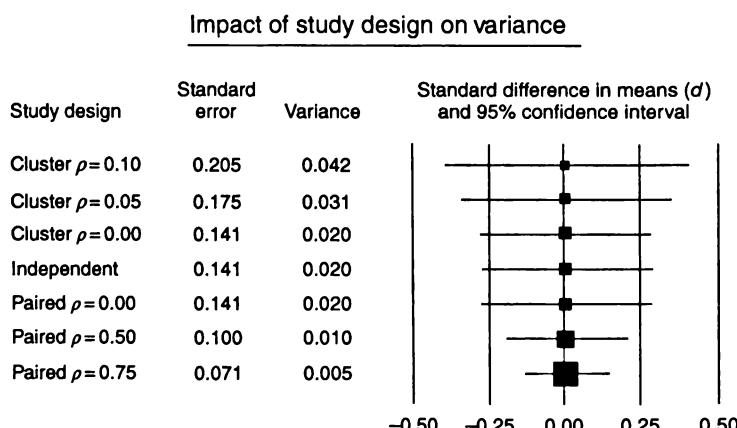
**Table 8.2** Impact of study design on variance.

Design	N per group	Intraclass correlation	Correlation	Standard error	Variance
A Cluster	10 × 10	0.10		0.205	0.042
B Cluster	10 × 10	0.05		0.175	0.031
C Cluster	10 × 10	0.00		0.141	0.020
D Independent	100			0.141	0.020
E Paired	100 pairs		0.00	0.141	0.020
F Paired	100 pairs		0.50	0.100	0.010
G Paired	100 pairs		0.75	0.071	0.005

are based on paired (or matched groups) with the same sample size (100 pairs). If the pre-post correlation was 0.00 (line E) then the matching would have no impact and the variance would remain at 0.02, but if the correlation was 0.50 (line F) or 0.75 (line G), then the variance would drop to 0.01 or 0.005.

Another design is the clustered trial, where an entire cluster of participants is assigned to one condition or another. For example, the design might call for students within classrooms, where an entire classroom is assigned to a single condition. Just as the use of matched pairs served to *decrease* the error term, the use of clusters serves to *increase* the error term, and a study that used clustered groups would typically have a larger variance than one with two independent groups. In clustered trials the intraclass correlation reflects the difference between clusters. If the intraclass correlation was 0.0 (line C) then the clustering would have no impact and the variance would remain at 0.02, but if the intraclass correlation was 0.05 (line B) or 0.10 (line A) the variance would increase to 0.03 or 0.04 (assuming 10 clusters of 10 subjects per group).

Again, the same information is presented graphically in Figure 8.2 where the larger blocks (and narrower confidence intervals) represent studies with more precise estimates.

**Figure 8.2** Impact of study design on variance.

## Concluding remarks

The information conveyed by precision is critically important in both primary studies and meta-analysis.

When we are working with individual studies the precision defines a range of likely values for the true effect. The precision, usually reported as a standard error or confidence interval, tells us how much confidence we can have in the effect size. To report that the effect size is 0.50 plus/minus 0.10 is very different than to report an effect size of 0.50 plus/minus 0.50.

As we turn our attention from the single study to the synthesis, our perspective shifts somewhat. A person performing a narrative review might look at a very precise study and decide to assign that study substantial weight in the analysis. This is formalized in the meta-analysis, with more weight being assigned to the more precise studies, as discussed in Part 4.

### SUMMARY POINTS

- The precision with which we estimate an effect size can be expressed as a standard error or confidence interval (in the same metric as the effect size itself) or as a variance (in a squared metric).
- The precision is driven primarily by the sample size, with larger studies yielding more precise estimates of the effect size.
- Other factors affecting precision include the study design, with matched groups yielding more precise estimates (as compared with independent groups) and clustered groups yielding less precise estimates.
- In addition to these general factors, there are unique factors that affect the precision for each effect size index.
- Studies that yield more precise estimates of the effect size carry more information and are assigned more weight in the meta-analysis.



# Concluding Remarks

While many meta-analyses use one of the effect sizes presented above, other options exist. Researchers working in medicine sometimes use the hazard ratio (based on the time to event in two groups) or the rate ratio (based on events by time in two groups). Nor are we limited to indices that look at the impact of a treatment or the relationship between two variables. Some indices simply report the mean, risk, or rate in a single group. For example, we could perform a meta-analysis of studies that had estimated the prevalence of HIV infection in different countries.

As we move on to formulas for meta-analysis we will be using one or another effect size as an example in each chapter. However, it is important to understand that once we have computed an effect size and variance for each study, the formulas for computing a summary effect, for assessing heterogeneity, and so on, are the same regardless of whether the effect size is a raw difference in means, a standardized difference in means, a log risk ratio, or another index.

### Further Reading

- Borenstein, M., Hedges L.V., Higgins, J.P.T., & Rothstein, H. (in preparation). *Computing Effect Sizes for Meta-analysis*. Chichester, UK: John Wiley & Sons, Ltd.\*
- Cooper, H., Hedges, L.V., & Valentine, J. (2019). *The Handbook of Research Synthesis*, 3rd edn. New York: Russell Sage Foundation.
- Deeks, J.J. (2002). Issues in the selection of a summary statistic for meta-analysis of clinical trials with binary outcomes. *Statistics in Medicine* 21: 1575–1600.
- Glass, G., McGaw, B., & Smith, M. (1981). *Meta-analysis in Social Research*. Newbury Park, CA: Sage.
- Hedges, L.V., Gurevitch, J., & Curtis, P. (1999). The meta-analysis of response ratios in experimental ecology. *Ecology* 80: 1150–1156.
- Higgins J.P.T., Thomas J., Chandler J., Cumpston M., Li T., Page M.J., Welch V.A. (editors) (2019). *Cochrane Handbook for Systematic Reviews of Interventions*. 2nd Edition. Chichester, UK: John Wiley & Sons, Ltd.
- Lipsey, M. & Wilson, D. (2001). *Practical Meta-analysis*. Thousand Oaks, CA: Sage.
- Rosenthal, R., Rosnow, R., & Rubin, D. (2000). *Contrasts and Effect Sizes in Behavioral Research: A Correlational Approach*. Cambridge, UK: Cambridge University Press.
- Shadish, W. (2003). *Effect Size Calculator*. St. Paul, MN: Assessment Systems Corporation.

---

\* Note. The first of these references (Borenstein *et al.*, in preparation) is the companion volume to this text, dedicated entirely to the computation of effect sizes and their variance.



# Fixed-Effect Versus Random-Effects Models



# Overview

---

Introduction  
Nomenclature

---

## INTRODUCTION

Most meta-analyses are based on one of two statistical models, the fixed-effect model or the random-effects model.

Under the fixed-effect model we assume that there is one *true effect size* (hence the term *fixed effect*) that underlies all the studies in the analysis, and that all differences in observed effects are due to sampling error. While we follow the practice of calling this a fixed-effect model, a more descriptive term would be a *common-effect* model. In either case, we use the singular (*effect*) since there is only one true effect.

By contrast, under the random-effects model we allow that the true effect could vary from study to study. For example, the effect size might be higher (or lower) in studies where the participants are older, or more educated, or healthier than in others, or when a more intensive variant of an intervention is used, and so on. Because studies will differ in the mixes of participants and in the implementations of interventions, among other reasons, there may be *different effect sizes* underlying different studies. If it were possible to perform an infinite number of studies (based on the inclusion criteria for our analysis), the true effect sizes for these studies would be distributed about some mean. The effect sizes in the studies that actually *were performed* are assumed to represent a random sample of these effect sizes (hence the term *random effects*). Here, we use the plural (*effects*) as there is an array of true effects.

In the chapters that follow we discuss the two models and show how to compute a summary effect using each one. Because the computations for a summary effect are not always intuitive, it helps to keep in mind that the summary effect is nothing more than the mean of the effect sizes, with more weight assigned to the more precise studies. We need to consider what we mean by the *more precise* studies and how this translates into a study weight (this depends on the model), but not lose track of the fact that we are simply computing a weighted mean.

## NOMENCLATURE

Throughout this Part we distinguish between a true effect size and an observed effect size. A study's *true effect size* is the effect size in the underlying population, and is the effect size that we would observe if the study had an infinitely large sample size (and therefore no sampling error). A study's *observed effect size* is the effect size that is actually observed.

In the schematics we use different symbols to distinguish between true effects and observed effects. For individual studies we use a circle for the former and a square for the latter (see Figure 10.1). For summary effects we use a triangle for the former and a diamond for the latter.

	True effect	Observed effect
Study	●	■
Combined	▼	◆

Figure 10.1 Symbols for true and observed effects.

## Worked examples

In meta-analysis the same formulas apply regardless of the effect-size index being used. To allow the reader to work with an effect size of their choosing, we have separated the formulas (which are presented in the following chapters) from the worked examples (which are presented in Chapter 14). There, we provide a worked example for the standardized mean difference, one for the odds ratio, and one for correlations.

The reader is encouraged to select one of the worked examples and follow the details of the computations while studying the formulas. The three datasets and all computations are available as Excel spreadsheets on the book's website

# Fixed-Effect Model

---

### Introduction

The true effect size

Impact of sampling error

Performing a fixed-effect meta-analysis

---

## INTRODUCTION

In this chapter we introduce the fixed-effect model. We discuss the assumptions of this model, and show how these are reflected in the formulas used to compute a summary effect, and in the meaning of the summary effect.

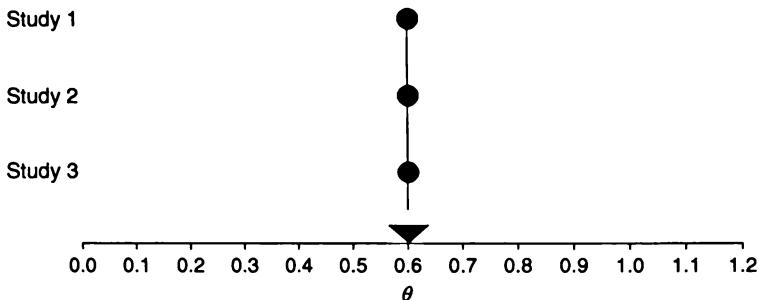
## THE TRUE EFFECT SIZE

Under the fixed-effect model we assume that all studies in the meta-analysis share a common (true) effect size. Put another way, all factors that could influence the effect size are the same in all the studies, and therefore the true effect size is the same (hence the label *fixed*) in all the studies. We denote the true (unknown) effect size by theta ( $\theta$ ).

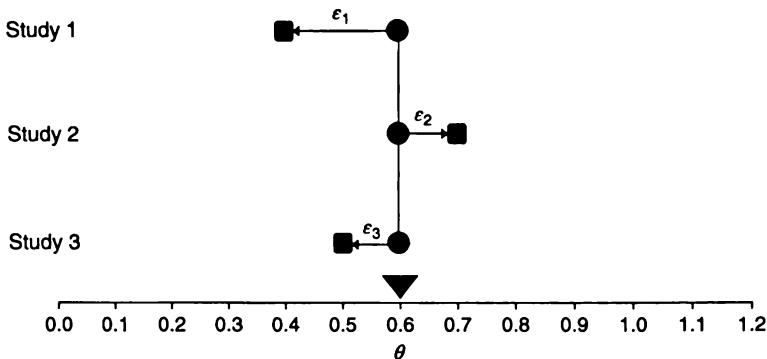
In Figure 11.1 the true overall effect size is 0.60 and this effect (represented by a triangle) is shown at the bottom. The true effect for each study is represented by a circle. Under the definition of a fixed-effect model the true effect size for each study must also be 0.60, and so these circles are aligned directly above the triangle.

## IMPACT OF SAMPLING ERROR

Since all studies share the same true effect, it follows that the observed effect size varies from one study to the next only because of the random error inherent in each study. If each study had an infinite sample size the sampling error would be zero and the observed effect for each study would be the same as the true effect. If we were to plot the observed effects rather than the true effects, the observed effects would exactly coincide with the true effects.



**Figure 11.1** Fixed-effect model – true effects.



**Figure 11.2** Fixed-effect model – true effects and sampling error.

In practice, of course, the sample size in each study is not infinite, and so there is sampling error and the effect observed in the study is not the same as the true effect. In Figure 11.2 the true effect for each study is still 0.60 (as depicted by the circles) but the observed effect (depicted by the squares) differs from one study to the next.

In Study 1 the sampling error ( $\epsilon_1$ ) is  $-0.20$ , which yields an observed effect ( $Y_1$ ) of

$$Y_1 = 0.60 - 0.20 = 0.40.$$

In Study 2 the sampling error ( $\epsilon_2$ ) is  $0.10$ , which yields an observed effect ( $Y_2$ ) of

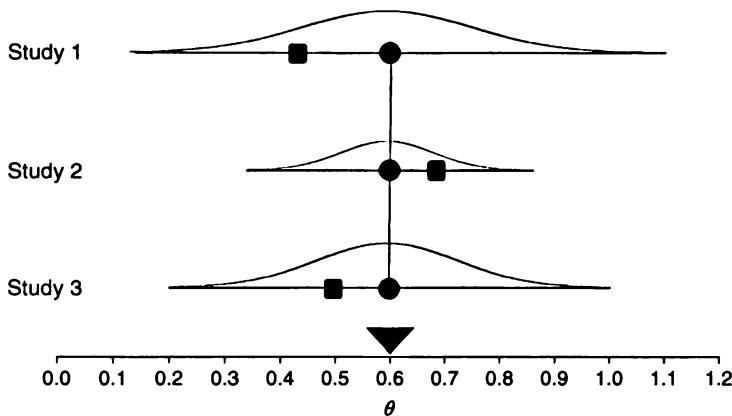
$$Y_2 = 0.60 + 0.10 = 0.70.$$

In Study 3 the sampling error ( $\epsilon_3$ ) is  $-0.10$ , which yields an observed effect ( $Y_3$ ) of

$$Y_3 = 0.60 - 0.10 = 0.50.$$

More generally, the observed effect  $Y_i$  for any study is given by the population mean plus the sampling error in that study. That is,

$$Y_i = \theta + \epsilon_i. \quad (11.1)$$



**Figure 11.3** Fixed-effect model – distribution of sampling error.

While the error in any given study is random, we *can* estimate the sampling distribution of the errors. In Figure 11.3 we have placed a normal curve about the true effect size for each study, with the width of the curve being based on the variance in that study. In Study 1 the sample size was small, the variance large, and the observed effect is likely to fall anywhere in the relatively wide range of 0.20 to 1.00. By contrast, in Study 2 the sample size was relatively large, the variance is small, and the observed effect is likely to fall in the relatively narrow range of 0.40 to 0.80. (The width of the normal curve is based on the square root of the variance, or standard error).

## PERFORMING A FIXED-EFFECT META-ANALYSIS

In an actual meta-analysis, of course, rather than starting with the population effect and making projections about the observed effects, we work backwards, starting with the observed effects and trying to estimate the population effect. In order to obtain the most precise estimate of the population effect (to minimize the variance) we compute a weighted mean, where the weight assigned to each study is the inverse of that study's variance. Concretely, the weight assigned to each study in a fixed-effect meta-analysis is

$$W_i = \frac{1}{V_{Y_i}}, \quad (11.2)$$

where  $V_{Y_i}$  is the within-study variance for study ( $i$ ). The weighted mean ( $M$ ) is then computed as

$$M = \frac{\sum_{i=1}^k W_i Y_i}{\sum_{i=1}^k W_i}, \quad (11.3)$$

that is, the sum of the products  $W_i Y_i$  (effect size multiplied by weight) divided by the sum of the weights.

The variance of the summary effect is estimated as the reciprocal of the sum of the weights, or

$$V_M = \frac{1}{\sum_{i=1}^k W_i}, \quad (11.4)$$

and the estimated standard error of the summary effect is then the square root of the variance,

$$SE_M = \sqrt{V_M}. \quad (11.5)$$

Then, 95% lower and upper limits for the summary effect are estimated as

$$LL_M = M - 1.96 \times SE_M \quad (11.6)$$

and

$$UL_M = M + 1.96 \times SE_M. \quad (11.7)$$

Finally, a Z-value to test the null hypothesis that the common true effect  $\theta$  is zero can be computed using

$$Z = \frac{M}{SE_M}. \quad (11.8)$$

For a one-tailed test the p-value is given by

$$p = 1 - \Phi(\pm|Z|), \quad (11.9)$$

where we choose '+' if the difference is in the expected direction and '-' otherwise, and for a two-tailed test by

$$p = 2[1 - (\Phi(|Z|))], \quad (11.10)$$

where  $\Phi(Z)$  is the standard normal cumulative distribution. This function is tabled in many introductory statistics books, and is implemented in Excel as the function =NORMSDIST(Z).

### Illustrative example

We suggest that you turn to a worked example for the fixed-effect model before proceeding to the random-effects model. A worked example for the standardized mean difference (Hedges'  $g$ ) is on page 81, a worked example for the odds ratio is on page 85, and a worked example for correlations is on page 90.

#### SUMMARY POINTS

- Under the fixed-effect model all studies in the analysis share a common true effect.
- The summary effect is our estimate of this common effect size, and the null hypothesis is that this common effect is zero (for a difference) or one (for a ratio).
- All observed dispersion reflects sampling error, and study weights are assigned with the goal of minimizing this within-study error.

# Random-Effects Model

---

Introduction

The true effect sizes

Impact of sampling error

Performing a random-effects meta-analysis

---

## INTRODUCTION

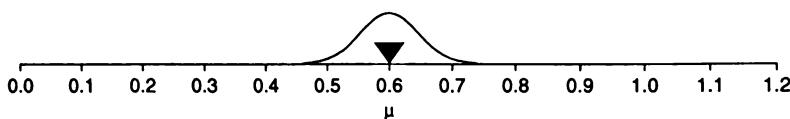
In this chapter we introduce the random-effects model. We discuss the assumptions of this model, and show how these are reflected in the formulas used to compute a summary effect, and in the meaning of the summary effect.

## THE TRUE EFFECT SIZES

The fixed-effect model, discussed above, starts with the assumption that the true effect size is the same in all studies. However, in many systematic reviews this assumption is implausible. When we decide to incorporate a group of studies in a meta-analysis, we assume that the studies have enough in common that it makes sense to synthesize the information, but there is generally no reason to assume that they are *identical* in the sense that the true effect size is *exactly the same* in all the studies.

For example, suppose that we are working with studies that compare the proportion of patients developing a disease in two groups (vaccinated versus a placebo). If the treatment works we would expect the effect size (say, the risk ratio) to be *similar but not identical* across studies. The effect size might be higher (or lower) when the participants are older, or more educated, or healthier than others, or when a more intensive variant of an intervention is used, and so on. Because studies will differ in the mixes of participants and in the implementations of interventions, among other reasons, there may be *different effect sizes* underlying different studies.

Or, suppose that we are working with studies that assess the impact of an educational intervention. The magnitude of the impact might vary depending on the other resources



**Figure 12.1** Random-effects model – distribution of the true effects.

available to the children, the class size, the age, and other factors, which are likely to vary from study to study.

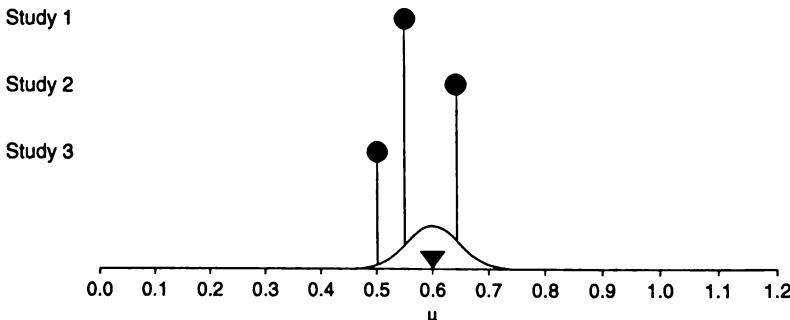
We might not have assessed these covariates in each study. Indeed, we might not even know what covariates actually are related to the size of the effect. Nevertheless, logic dictates that such factors do exist and will lead to variations in the magnitude of the effect.

One way to address this variation across studies is to perform a *random-effects* meta-analysis. In a random-effects meta-analysis we usually assume that the true effects are normally distributed. For example, in Figure 12.1 the mean of all true effect sizes is 0.60 but the individual effect sizes are distributed about this mean, as indicated by the normal curve. The width of the curve suggests that most of the true effects fall in the range of 0.50 to 0.70.

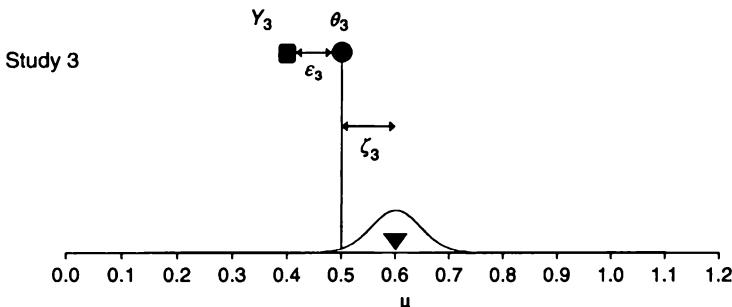
## IMPACT OF SAMPLING ERROR

Suppose that our meta-analysis includes three studies drawn from the distribution of studies depicted by the normal curve, and that the true effects (denoted  $\theta_1$ ,  $\theta_2$ , and  $\theta_3$ ) in these studies happen to be 0.50, 0.55 and 0.65 (see Figure 12.2).

If each study had an infinite sample size the sampling error would be zero and the observed effect for each study would be the same as the true effect for that study. If we were to plot the observed effects rather than the true effects, the observed effects would exactly coincide with the true effects.



**Figure 12.2** Random-effects model – true effects.



**Figure 12.3** Random-effects model – true and observed effect in one study.

Of course, the sample size in any study is not infinite and therefore the sampling error is not zero. If the true effect size for a study is  $\theta_i$ , then the observed effect for that study will be less than or greater than  $\theta_i$  because of sampling error. For example, consider Study 3 in Figure 12.2. This study is the subject of Figure 12.3, where we consider the factors that control the observed effect. The true effect for Study 3 is 0.50 but the sampling error for this study is -0.10, and the observed effect for this study is 0.40.

This figure also highlights the fact that the distance between the overall mean and the observed effect in any given study consists of two distinct parts: true variation in effect sizes ( $\zeta_i$ ) and sampling error ( $\epsilon_i$ ). In Study 3 the total distance from  $\mu$  to  $\theta_3$  is -0.20. The distance from  $\mu$  to  $\theta_3$  (0.60 to 0.50) reflects the fact that the true effect size actually varies from one study to the next, while the distance from  $\theta_3$  to  $Y_3$  (0.5 to 0.4) is sampling error.

More generally, the observed effect  $Y_i$  for any study is given by the grand mean, the deviation of the study's true effect from the grand mean, and the deviation of the study's observed effect from the study's true effect. That is,

$$Y_i = \mu + \zeta_i + \epsilon_i \quad (12.1)$$

Therefore, to predict how far the observed effect  $Y_i$  is likely to fall from  $\mu$  in any given study we need to consider both the variance of  $\zeta_i$ , and the variance of  $\epsilon_i$ .

The distance from  $\mu$  (the triangle) to each  $\theta_i$  (the circles) depends on the standard deviation of the distribution of the true effects across studies, called  $\tau$  (tau) (or  $\tau^2$  for its variance). The same value of  $\tau^2$  applies to all studies in the meta-analysis, and in Figure 12.4 is represented by the normal curve at the bottom, which extends roughly from 0.50 to 0.70.

The distance from  $\theta_i$  to  $Y_i$ , depends on the sampling distribution of the sample effects about  $\theta_i$ . This depends on the variance of the observed effect size from each study,  $V_{Y_i}$ , and so will vary from one study to the next. In Figure 12.4 the curve for Study 1 is relatively wide while the curve for Study 2 is relatively narrow.

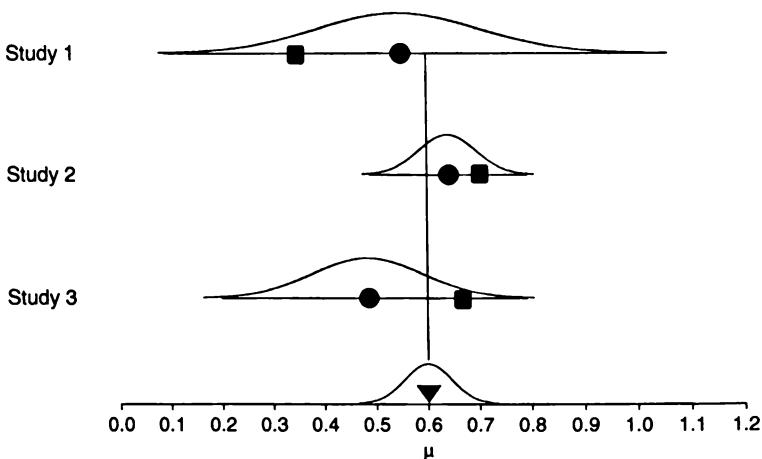


Figure 12.4 Random-effects model – between-study and within-study variance.

## PERFORMING A RANDOM-EFFECTS META-ANALYSIS

In an actual meta-analysis, of course, rather than start with the population effect and make projections about the observed effects, we start with the observed effects and try to estimate the population effect. In other words our goal is to use the collection of  $Y_i$  to estimate the overall mean,  $\mu$ . In order to obtain the most precise estimate of the overall mean (to minimize the variance) we compute a weighted mean, where the weight assigned to each study is the inverse of that study's variance.

To compute a study's variance under the random-effects model, we need to know both the within-study variance and  $\tau^2$ , since the study's total variance is the sum of these two values. Formulas for computing the within-study variance were presented in Part 3. A method for estimating the between-studies variance is given here so that we can proceed with the worked example, but a full discussion of this method is deferred to Part 4, where we shall pursue the issue of heterogeneity in some detail.

### Estimating tau-squared

The parameter  $\tau^2$  (tau-squared) is the between-studies variance (the variance of the effect size parameters across the population of studies). In other words, if we somehow knew the *true* effect size for each study, and computed the variance of these effect sizes (across an infinite number of studies), this variance would be  $\tau^2$ . One method for estimating  $\tau^2$  is the method of moments (or the DerSimonian and Laird) method, as follows. We compute

$$\tau^2 = \frac{Q - df}{C}, \quad (12.2)$$

where

$$Q = \sum_{i=1}^k W_i Y_i^2 - \frac{\left( \sum_{i=1}^k W_i Y_i \right)^2}{\sum_{i=1}^k W_i}, \quad (12.3)$$

$$df = k - 1, \quad (12.4)$$

where  $k$  is the number of studies, and

$$C = \sum W_i - \frac{\sum W_i^2}{\sum W_i}. \quad (12.5)$$

### Estimating the mean effect size

In the fixed-effect analysis each study was weighted by the inverse of its variance. In the random-effects analysis, too, each study will be weighted by the inverse of its variance. The difference is that the variance now includes the original (within-studies) variance plus the estimate of the between-studies variance,  $\tau^2$ . In keeping with the book's convention, we use  $\tau^2$  to refer to the parameter and  $T^2$  to refer to the sample estimate of that parameter.

To highlight the parallel between the formulas here (random effects) and those in the previous chapter (fixed effect) we use the same notations but add an asterisk (\*) to represent the random-effects version. Under the random-effects model the weight assigned to each study is

$$W_i^* = \frac{1}{V_{Y_i}^*}, \quad (12.6)$$

where  $V_{Y_i}^*$  is the within-study variance for study  $i$  plus the between-studies variance,  $T^2$ . That is,

$$V_{Y_i}^* = V_{Y_i} + T^2.$$

The weighted mean,  $M^*$ , is then computed as

$$M^* = \frac{\sum_{i=1}^k W_i^* Y_i}{\sum_{i=1}^k W_i^*}. \quad (12.7)$$

that is, the sum of the products (effect size multiplied by weight) divided by the sum of the weights.

The variance of the summary effect is estimated as the reciprocal of the sum of the weights, or

$$V_{M^*} = \frac{1}{\sum_{i=1}^k W_i^*}, \quad (12.8)$$

and the estimated standard error of the summary effect is then the square root of the variance,

$$SE_{M^*} = \sqrt{V_{M^*}}. \quad (12.9)$$

The 95% lower and upper limits for the summary effect would be computed as

$$LL_{M^*} = M^* - 1.96 \times SE_{M^*}, \quad (12.10)$$

and

$$UL_{M^*} = M^* + 1.96 \times SE_{M^*}. \quad (12.11)$$

Finally, a Z-value to test the null hypothesis that the mean effect  $\mu$  is zero could be computed using

$$Z^* = \frac{M^*}{SE_{M^*}}. \quad (12.12)$$

For a one-tailed test the  $p$ -value is given by

$$p^* = 1 - \Phi(\pm|Z^*|), \quad (12.13)$$

where we choose '+' if the difference is in the expected direction or '-' otherwise, and for a two-tailed test by

$$p^* = 2[1 - (\Phi(|Z^*|))], \quad (12.14)$$

where  $\Phi(Z^*)$  is the standard normal cumulative distribution. This function is tabled in many introductory statistics books, and is implemented in Excel as the function =NORMSDIST(Z\*).

### **Illustrative example**

As before, we suggest that you turn to one of the worked examples in the next chapter before proceeding with this discussion.

### **SUMMARY POINTS**

- Under the random-effects model, the true effects in the studies are assumed to have been sampled from a distribution of true effects.
- The summary effect is our estimate of the mean of all relevant true effects, and the null hypothesis is that the mean of these effects is 0.0 (equivalent to a ratio of 1.0 for ratio measures).
- Since our goal is to estimate the mean of the distribution, we need to take account of two sources of variance. First, there is within-study error in estimating the effect in each study. Second (even if we knew the true mean for each of our studies), there is variation in the true effects across studies. Study weights are assigned with the goal of minimizing both sources of variance.

# Fixed-Effect Versus Random-Effects Models

---

### Introduction

Definition of a summary effect

Estimating the summary effect

Extreme effect size in a large study or a small study

Confidence interval

The null hypothesis

Which model should we use?

Model should not be based on the test for heterogeneity

Concluding remarks

---

## INTRODUCTION

In Chapter 11 and Chapter 12 we introduced the fixed-effect and random-effects models. Here, we highlight the conceptual and practical differences between them.

Consider the forest plots in Figures 13.1 and 13.2. They include the same six studies, but the first uses a fixed-effect analysis and the second a random-effects analysis. These plots provide a context for the discussion that follows.

## DEFINITION OF A SUMMARY EFFECT

Both plots show a summary effect on the bottom line, but the meaning of this summary effect is different in the two models. In the fixed-effect analysis we assume that the true effect size is the same in all studies, and the summary effect is our estimate of this common effect size. In the random-effects analysis we assume that the true effect size varies from one study to the next, and that the studies in our analysis represent a random sample of effect sizes that could have been observed. The summary effect is our estimate of the mean of these effects.

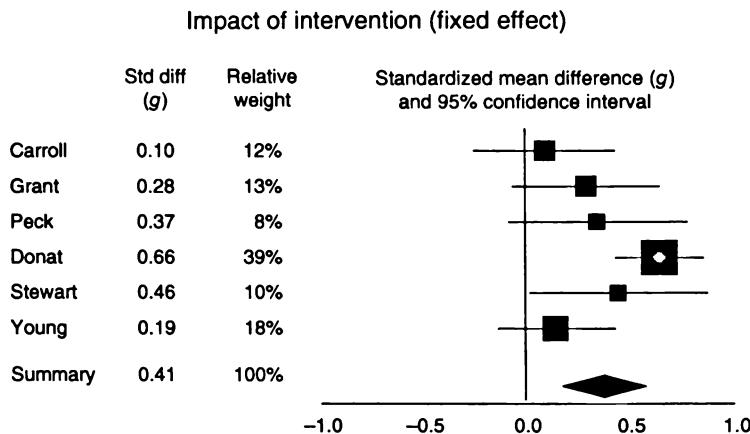


Figure 13.1 Fixed-effect model – forest plot showing relative weights.

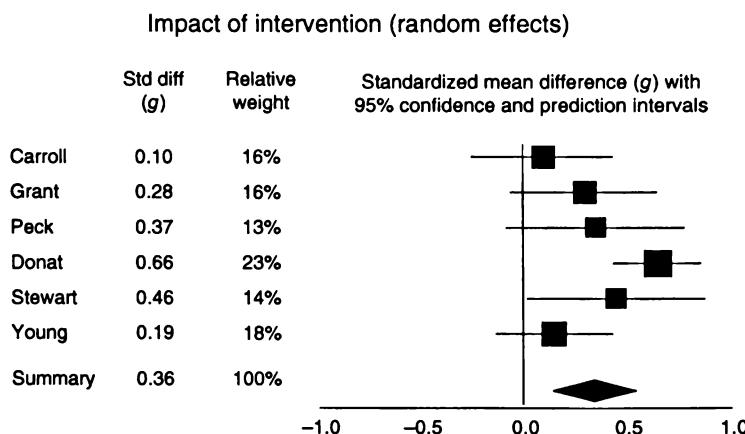


Figure 13.2 Random-effects model – forest plot showing relative weights.

## ESTIMATING THE SUMMARY EFFECT

Under the fixed-effect model we assume that the true effect size for all studies is identical, and the only reason the effect size varies between studies is sampling error (error in estimating the effect size). Therefore, when assigning weights to the different studies we can largely ignore the information in the smaller studies since we have better information about the same effect size in the larger studies.

By contrast, under the random-effects model the goal is not to estimate one true effect, but to estimate the mean of a distribution of effects. Since each study provides information about a different effect size, we want to be sure that all these effect sizes are represented in the summary estimate. This means that we cannot discount a small

study by giving it a very small weight (the way we would in a fixed-effect analysis). The estimate provided by that study may be imprecise, but it is information about an effect that no other study has estimated. By the same logic we cannot give too much weight to a very large study (the way we might in a fixed-effect analysis). Our goal is to estimate the mean effect in a range of studies, and we do not want that overall estimate to be overly influenced by any one of them.

In these graphs, the weight assigned to each study is reflected in the size of the box (specifically, the area) for that study. Under the fixed-effect model there is a wide range of weights (as reflected in the size of the boxes) whereas under the random-effects model the weights fall in a relatively narrow range. For example, compare the weight assigned to the largest study (Donat) with that assigned to the smallest study (Peck) under the two models. Under the fixed-effect model Donat is given about five times as much weight as Peck. Under the random-effects model Donat is given only 1.8 times as much weight as Peck.

### EXTREME EFFECT SIZE IN A LARGE STUDY OR A SMALL STUDY

How will the selection of a model influence the overall effect size? In this example Donat is the largest study, and also happens to have the highest effect size. Under the fixed-effect model Donat was assigned a large share (39%) of the total weight and pulled the mean effect up to 0.41. By contrast, under the random-effects model Donat was assigned a relatively modest share of the weight (23%). It therefore had less pull on the mean, which was computed as 0.36.

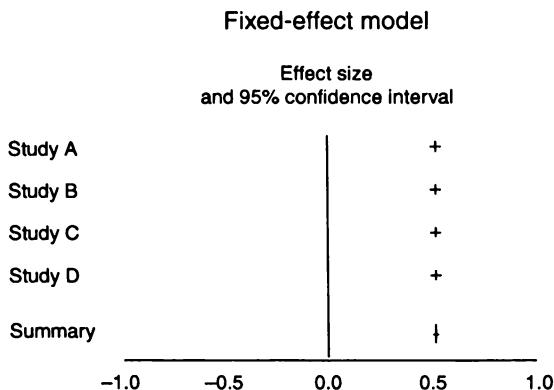
Similarly, Carroll is one of the smaller studies and happens to have the smallest effect size. Under the fixed-effect model Carroll was assigned a relatively small proportion of the total weight (12%), and had little influence on the summary effect. By contrast, under the random-effects model Carroll carried a somewhat higher proportion of the total weight (16%) and was able to pull the weighted mean toward the left.

The operating premise, as illustrated in these examples, is that whenever  $\tau^2$  is nonzero, the relative weights assigned under random effects will be *more balanced* than those assigned under fixed effects. As we move from fixed effect to random effects, extreme studies will lose influence if they are large, and will gain influence if they are small.

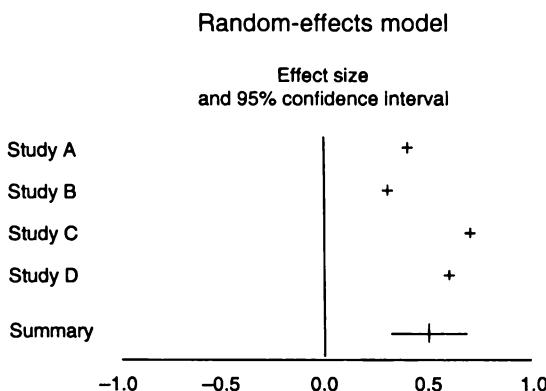
### CONFIDENCE INTERVAL

Under the fixed-effect model the only source of uncertainty is the within-study (sampling or estimation) error. Under the random-effects model there is this same source of uncertainty plus an additional source (between-studies variance). It follows that the variance, standard error, and confidence interval for the summary effect will always be larger (or wider) under the random-effects model than under the fixed-effect model (unless  $\tau^2$  is zero, in which case the two models are the same). In this example, the standard error is 0.064 for the fixed-effect model, and 0.105 for the random-effects model.

Consider what would happen if we had five studies, and each study had an infinitely large sample size. Under either model the confidence interval for the effect size in each



**Figure 13.3** Very large studies under fixed-effect model.



**Figure 13.4** Very large studies under random-effects model.

study would have a width approaching zero, since we know the effect size in that study with perfect precision. Under the fixed-effect model the summary effect would also have a confidence interval with a width of zero, since we know the common effect precisely (Figure 13.3). By contrast, under the random-effects model the width of the confidence interval would not approach zero (Figure 13.4). While we know the effect in each study precisely, these effects have been sampled from a universe of possible effect sizes, and provide only an estimate of the mean effect. Just as the error within a study will approach zero only as the sample size approaches infinity, so too the error of these studies as an estimate of the mean effect will approach zero only as the number of studies approaches infinity.

More generally, it is instructive to consider what factors influence the standard error of the summary effect under the two models. The following formulas are based on a meta-analysis of means from  $k$  one-group studies, but the conceptual argument applies

to all meta-analyses. The within-study variance of each mean depends on the standard deviation (denoted  $\sigma$ ) of participants' scores and the sample size of each study ( $n$ ). For simplicity we assume that all of the studies have the same sample size and the same standard deviation (see Box 13.1 for details).

Under the fixed-effect model the standard error of the summary effect is given by

$$SE_M = \sqrt{\frac{\sigma^2}{k \times n}} \quad (13.1)$$

It follows that with a large enough sample size the standard error will approach zero, and this is true whether the sample size is concentrated on one or two studies, or dispersed across any number of studies.

Under the random-effects model the standard error of the summary effect is given by

$$SE_M = \sqrt{\frac{\sigma^2}{k \times n} + \frac{\tau^2}{k}} \quad (13.2)$$

The first term is identical to that for the fixed-effect model and, again, with a large enough sample size, this term will approach zero. By contrast, the second term (which reflects the between-studies variance) will only approach zero as the number of studies approaches infinity. These formulas do not apply exactly in practice, but the conceptual argument does. Namely, increasing the sample size within studies is not sufficient to reduce the standard error beyond a certain point (where that point is determined by  $\tau^2$  and  $k$ ). If there is only a small number of studies, then the standard error could still be substantial even if the total  $n$  is in the tens of thousands or higher.

### BOX 13.1 FACTORS THAT INFLUENCE THE STANDARD ERROR OF THE SUMMARY EFFECT

To illustrate the concepts with some simple formulas, let us consider a meta-analysis of studies with the very simplest design, such that each study comprises a single sample of  $n$  observations with standard deviation  $\sigma$ . We combine estimates of the mean in a meta-analysis. The variance of each estimate is

$$V_{Y_i} = \frac{\sigma^2}{n}$$

so the (inverse-variance) weight in a fixed-effect meta-analysis is

$$W_i = \frac{1}{\sigma^2/n} = \frac{n}{\sigma^2}$$

and the variance of the summary effect under the fixed-effect model the standard error is given by

$$V_M = \frac{1}{\sum_{i=1}^k W_i} = \frac{1}{k \times n/\sigma^2} = \frac{\sigma^2}{k \times n}$$

**BOX 13.1 CONTINUED**

Therefore under the fixed-effect model the (true) standard error of the summary mean is given by

$$SE_M = \sqrt{\frac{\sigma^2}{k \times n}}$$

Under the random-effects model the weight awarded to each study is

$$W_i^* = \frac{1}{(\sigma^2/n) + \tau^2}$$

and the (true) standard error of the summary mean turns out to be

$$SE_M^* = \sqrt{\frac{\sigma^2}{k \times n} + \frac{\tau^2}{k}}.$$

**THE NULL HYPOTHESIS**

Often, after computing a summary effect, researchers perform a test of the null hypothesis. Under the fixed-effect model the null hypothesis being tested is that there is zero effect in *every study*. Under the random-effects model the null hypothesis being tested is that the *mean effect* is zero. Although some may treat these hypotheses as interchangeable, they are in fact different, and it is imperative to choose the test that is appropriate to the inference a researcher wishes to make.

**WHICH MODEL SHOULD WE USE?**

The selection of a computational model should be based on our expectation about whether or not the studies share a common effect size and on our goals in performing the analysis.

**Fixed effect**

It makes sense to use the fixed-effect model if two conditions are met. First, we believe that all the studies included in the analysis are functionally identical. Second, our goal is to compute the common effect size for the identified population, and not to generalize to other populations.

For example, suppose that a pharmaceutical company will use a thousand patients to compare a drug versus a placebo. Because the staff can work with only 100 patients at a time, the company will run a series of ten trials with 100 patients in each. The studies are identical in the sense that any variables which can have an impact on the outcome are the same across the ten studies. Specifically, the studies draw patients from a common pool, using the same researchers, dose, measure, and so on (we assume that there is no concern about practice effects for the researchers, nor for the different

starting times of the various cohorts). All the studies are expected to share a common effect and so the first condition is met. The goal of the analysis is to see if the drug works in the population from which the patients were drawn (and not to extrapolate to other populations), and so the second condition is met, as well.

In this example the fixed-effect model is a plausible fit for the data and meets the goal of the researchers. It should be clear, however, that this situation is relatively rare. The vast majority of cases will more closely resemble those discussed immediately below.

### Random effects

By contrast, when the researcher is accumulating data from a series of studies that had been performed by researchers operating independently, it would be unlikely that all the studies were functionally equivalent. Typically, the subjects or interventions in these studies would have differed in ways that would have impacted on the results, and therefore we should not assume a common effect size. Therefore, in these cases the random-effects model is more easily justified than the fixed-effect model.

Additionally, the goal of this analysis is usually to generalize to a range of scenarios. Therefore, if one did make the argument that all the studies used an identical, narrowly defined population, then it would not be possible to extrapolate from this population to others, and the utility of the analysis would be severely limited.

### A caveat

There is one caveat to the above. If the number of studies is very small, then the estimate of the between-studies variance ( $\tau^2$ ) will have poor precision. While the random-effects model is still the appropriate model, we lack the information needed to apply it correctly. In this case the reviewer may choose among several options, each of them problematic.

One option is to report the separate effects and not report a summary effect. The hope is that the reader will understand that we cannot draw conclusions about the effect size and its confidence interval. The problem is that some readers will revert to vote counting (see Chapter 33) and possibly reach an erroneous conclusion.

Another option is to perform a fixed-effect analysis. This approach would yield a descriptive analysis of the included studies, but would not allow us to make inferences about a wider population. The problem with this approach is that (a) we do want to make inferences about a wider population and (b) readers will make these inferences even if they are not warranted.

A third option is to take a Bayesian approach, where the estimate of  $\tau^2$  is based on data from outside of the current set of studies. This is probably the best option, but the problem is that relatively few researchers have expertise in Bayesian meta-analysis. Additionally, some researchers have a philosophical objection to this approach.

For a more general discussion of this issue see *When does it make sense to perform a meta-analysis* in Chapter 45.

## MODEL SHOULD NOT BE BASED ON THE TEST FOR HETEROGENEITY

In the next chapter we will introduce a test of the null hypothesis that the between-studies variance is zero. This test is based on the amount of between-studies variance observed, relative to the amount we would expect if the studies actually shared a common effect size.

Some have adopted the practice of starting with a fixed-effect model and then switching to a random-effects model if the test of homogeneity is statistically significant. This practice should be strongly discouraged because the decision to use the random-effects model should be based on our understanding of whether or not all studies share a common effect size, and not on the outcome of a statistical test (especially since the test for heterogeneity often suffers from low power).

If the study effect sizes are seen as having been sampled from a *distribution* of effect sizes, then the random-effects model, which reflects this idea, is the logical one to use. If the between-studies variance is substantial (and statistically significant) then the fixed-effect model is inappropriate. However, even if the between-studies variance does not meet the criterion for statistical significance (which may be due simply to low power) we should still take account of this variance when assigning weights. If  $\tau^2$  turns out to be zero, then the random-effects analysis reduces to the fixed-effect analysis, and so there is no cost to using this model.

On the other hand, if one has elected to use the fixed-effect model *a priori* but the test of homogeneity is statistically significant, then it might be helpful to revisit the assumptions that led to the selection of a fixed-effect model.

Rice, Higgins, and Lumley (2018) have suggested that it would be helpful to distinguish between two versions of the fixed-effect analysis. The label ‘fixed-effect’ where ‘effect’ is in the singular applies to the case where all studies share a common true effect size and our goal is to make an inference about the one population represented in the analysis. The label ‘fixed-effects’ where ‘effects’ is in the plural applies to the case where the true effect size varies across studies and our goal is to make an inference to the set of populations included in the analysis. In both cases, our inference is limited to the population (or populations) included in the analysis, and the computational formula (and results) are identical under the two models. By contrast, under the random-effects model, we use the studies in the analysis to make an inference to a wider universe of comparable studies.

In this volume, we will focus on the fixed-effect (singular) model and the random-effects model. However, the reader should be aware that the fixed-effects (plural) model is also an option, and there are situations where we may prefer to use this model. This might be the case if we do not have a sufficient number of studies to employ the random-effects model reliably. This might also be the case if the analysis is being used to support the approval of a new drug, and we need to make an inference to the studies in the analysis, without generalizing to a larger universe. In these cases, researchers have sometimes tried to justify the use of the fixed-effect model by arguing that the studies share a common effect size, when in fact the studies

do not share a common effect size. A better option is to acknowledge that the effect size varies across studies and apply the fixed-effects (plural) model.

For a lengthy discussion of this issue, see Borenstein (2019); Rice et al. (2018).

## CONCLUDING REMARKS

Our discussion of differences between the fixed-effect model and the random-effects model focused largely on the computation of a summary effect and the confidence intervals for the summary effect. We did not address the implications of the dispersion itself. Under the fixed-effect model we assume that all dispersion in observed effects is due to sampling error, but under the random-effects model we allow that some of that dispersion reflects real differences in effect size across studies. In the chapters that follow we discuss methods to quantify that dispersion and to consider its substantive implications.

### SUMMARY POINTS

- A fixed-effect meta-analysis estimates a single effect that is assumed to be common to every study, while a random-effects meta-analysis estimates the mean of a distribution of effects.
- Study weights are more balanced under the random-effects model than under the fixed-effect model. Large studies are assigned less relative weight and small studies are assigned more relative weight as compared with the fixed-effect model.
- The standard error of the summary effect and (it follows) the confidence intervals for the summary effect are wider under the random-effects model than under the fixed-effect model.
- The selection of a model must be based solely on the question of which model fits the distribution of effect sizes, and takes account of the relevant source(s) of error. When studies are gathered from the published literature, the random-effects model is generally a more plausible match.
- The strategy of starting with a fixed-effect model and then moving to a random-effects model if the test for heterogeneity is significant is a mistake, and should be strongly discouraged.



# Worked Examples (Part 1)

---

### Introduction

Worked example for continuous data (Part 1)

Worked example for binary data (Part 1)

Worked example for correlational data (Part 1)

---

## INTRODUCTION

In this chapter we present worked examples for continuous data (using the standardized mean difference), binary data (using the odds ratio) and correlational data (using the Fisher's  $z$  transformation).

All of the data sets and all computations are available as Excel spreadsheets on the book's website ([www.Introduction-to-Meta-Analysis.com](http://www.Introduction-to-Meta-Analysis.com)).

## WORKED EXAMPLE FOR CONTINUOUS DATA (PART 1)

In this example we start with the mean, standard deviation, and sample size, and will use the bias-corrected standardized mean difference (Hedges'  $g$ ) as the effect size measure.

### Summary data

The summary data for six studies are presented in Table 14.1.

### Compute the effect size and its variance for each study

The first step is to compute the effect size ( $g$ ) and variance for each study using the formulas in Chapter 4 (see (4.18) to (4.24)). For the first study (Carroll) we compute the pooled within-groups standard deviation

$$S_{\text{within}} = \sqrt{\frac{(60 - 1) \times 22^2 + (60 - 1) \times 20^2}{60 + 60 - 2}} = 21.0238.$$

**Table 14.1** Dataset 1 – Part A (basic data).

Study	Treated			Control		
	Mean	SD	n	Mean	SD	n
Carroll	94	22	60	92	20	60
Grant	98	21	65	92	22	65
Peck	98	28	40	88	26	40
Donat	94	19	200	82	17	200
Stewart	98	21	50	88	22	45
Young	96	21	85	92	22	85

Then we compute the standardized mean difference,  $d$ , and its variance as

$$d_1 = \frac{94 - 92}{21.0238} = 0.0951,$$

and

$$V_{d_1} = \frac{60 + 60}{60 \times 60} + \frac{0.0951^2}{2(60 + 60)} = 0.0334.$$

The correction factor ( $J$ ) is estimated as

$$J = \left(1 - \frac{3}{4 \times 118 - 1}\right) = 0.9936.$$

Finally, the bias-corrected standardized mean difference, Hedges'  $g$ , and its variance are given by

$$g_1 = 0.9936 \times 0.0951 = 0.0945,$$

and

$$V_{g_1} = 0.9936^2 \times 0.0334 = 0.0329.$$

This procedure is repeated for all six studies.

### Compute the summary effect using the fixed-effect model

The effect size and its variance are copied into Table 14.2 where they are assigned the generic labels  $Y$  and  $V_Y$ . We then compute the other values shown in the table. For Carroll,

$$W_1 = \frac{1}{0.0329} = 30.3515,$$

$$W_1 Y_1 = 30.3515 \times 0.0945 = 2.8690,$$

and so on for the other five studies. The sum of  $W$  is 244.215 and the sum of  $WY$  is 101.171. From these numbers we can compute the summary effect and related statistics, using formulas from Part 3 as follows (see (11.3) to (11.10)). In the computations that follow we use the generic  $M$  to represent Hedges'  $g$ .

$$M = \frac{101.171}{244.215} = 0.4143,$$

**Table 14.2** Dataset 1 – Part B (fixed-effect computations).

Study	Effect size $\gamma$	Variance within $V_\gamma$	Weight $W$	Calculated quantities		
				$W\gamma$	$W\gamma^2$	$W^2$
Carroll	0.095	0.033	30.352	2.869	0.271	921.214
Grant	0.277	0.031	32.568	9.033	2.505	1060.682
Peck	0.367	0.050	20.048	7.349	2.694	401.931
Donat	0.664	0.011	95.111	63.190	41.983	9046.013
Stewart	0.462	0.043	23.439	10.824	4.999	549.370
Young	0.185	0.023	42.698	7.906	1.464	1823.115
Sum			244.215	101.171	53.915	13802.325

$$V_M = \frac{1}{244.215} = 0.0041,$$

$$SE_M = \sqrt{0.0041} = 0.0640,$$

$$LL_M = 0.4143 - 1.96 \times 0.0640 = 0.2889,$$

$$UL_M = 0.4143 + 1.96 \times 0.0640 = 0.5397,$$

and

$$Z = \frac{0.4143}{0.0640} = 6.4739.$$

For a one-tailed test the  $p$ -value is given by

$$p = 1 - \Phi(6.4739) < 0.0001,$$

and for a two-tailed test, by

$$p = 2[1 - \Phi(|6.4739|)] < 0.0001.$$

In words, using fixed-effect weights, the standardized mean difference (Hedges'  $g$ ) is 0.41 with a 95% confidence interval of 0.29 to 0.54. The  $Z$ -value is 6.47, and the  $p$ -value is  $<0.0001$  (one-tailed) or  $<0.0001$  (two tailed). These results are illustrated in Figure 14.1.

### Compute an estimate of $\tau^2$

To estimate  $\tau^2$ , the variance of the true standardized mean differences, we use the DerSimonian and Laird method (see (12.2) to (12.5)). Using sums from Table 14.2,

$$Q = 53.915 - \left( \frac{101.171^2}{244.215} \right) = 12.0033,$$

$$df = (6 - 1) = 5,$$

$$C = 244.215 - \left( \frac{13802.325}{244.215} \right) = 187.698,$$

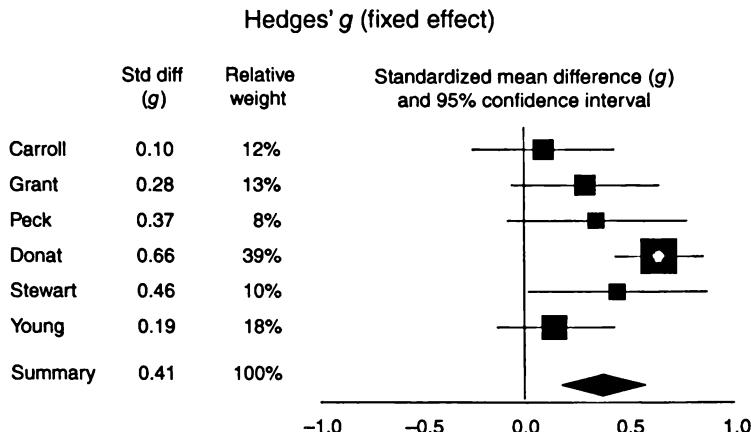


Figure 14.1 Forest plot of Dataset 1 – fixed-effect weights.

and

$$T^2 = \frac{12.0033 - 5}{187.698} = 0.0373.$$

### Compute the summary effect using the random-effects model

To compute the summary effect using the random-effects model we use the same formulas as for the fixed effect, but the variance for each study is now the sum of the variance within studies plus the variance between studies (see (12.6) to (12.13)).

For Carroll,

$$W_1^* = \frac{1}{(0.0329 + 0.0373)} = \frac{1}{(0.070)} = 14.2331,$$

and so on for the other studies as shown in Table 14.3. Note that the within-study variance is unique for each study, but there is only one value of  $\tau^2$ , so this value (estimated as 0.037) is applied to all studies.

Then,

$$M^* = \frac{32.342}{90.284} = 0.3582, \quad (14.1)$$

$$V_{M^*} = \frac{1}{90.284} = 0.0111, \quad (14.2)$$

$$SE_{M^*} = \sqrt{0.0111} = 0.1052,$$

$$LL_{M^*} = 0.3582 - 1.96 \times 0.1052 = 0.1520,$$

$$UL_{M^*} = 0.3582 + 1.96 \times 0.1052 = 0.5645,$$

$$Z^* = \frac{0.3582}{0.1052} = 3.4038,$$

**Table 14.3** Dataset 1 – Part C (random-effects computations).

Study	Effect size $Y$	Variance within $V_Y$	Variance between $T^2$	Variance total $V_Y + T^2$	Weight $W^*$	Calculated quantities $W^* Y$
Carroll	0.095	0.033	0.037	0.070	14.233	1.345
Grant	0.277	0.031	0.037	0.068	14.702	4.078
Peck	0.367	0.050	0.037	0.087	11.469	4.204
Donat	0.664	0.011	0.037	0.048	20.909	13.892
Stewart	0.462	0.043	0.037	0.080	12.504	5.774
Young	0.185	0.023	0.037	0.061	16.466	3.049
Sum					90.284	32.342

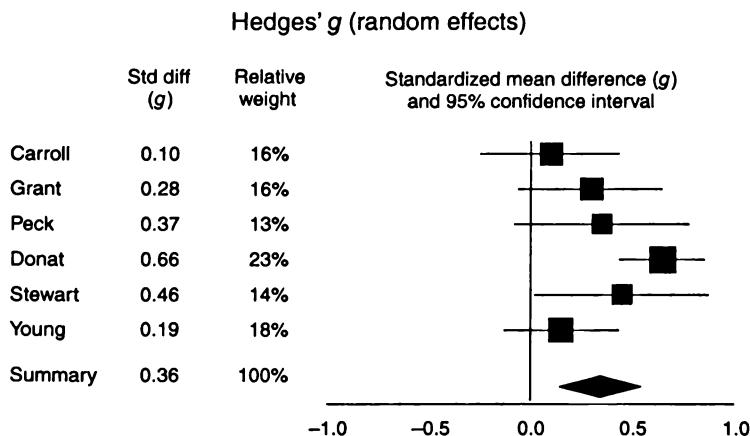
and, for a one-tailed test

$$p^* = 1 - \Phi(3.4038) = 0.0003$$

or, for a two-tailed test

$$p^* = 2[1 - \Phi(|3.4038|)] = 0.0007.$$

In words, using random-effect weights, the standardized mean difference (Hedges'  $g$ ) is 0.36 with a 95% confidence interval of 0.15 to 0.56. The Z-value is 3.40, and the  $p$ -value is 0.0003 (one-tailed) or 0.0007 (two-tailed). These results are illustrated in Figure 14.2.

**Figure 14.2** Forest plot of Dataset 1 – random-effects weights.

### WORKED EXAMPLE FOR BINARY DATA (PART 1)

In this example we start with the events and non-events in two independent groups and will use the odds ratio as the effect size measure.

**Table 14.4** Dataset 2 – Part A (basic data).

Study	Treated			Control		
	Events	Non-events	n	Events	Non-events	n
Saint	12	53	65	16	49	65
Kelly	8	32	40	10	30	40
Pilbeam	14	66	80	19	61	80
Lane	25	375	400	80	320	400
Wright	8	32	40	11	29	40
Day	16	49	65	18	47	65

### Summary data

The summary data for six studies is presented in Table 14.4.

### Compute the effect size and its variance for each study

For an odds ratio all computations are carried out using the log transformed values (see formulas (5.8) to (5.10)). For the first study (Saint) we compute the odds ratio, then the log odds ratio and its variance as

$$OddsRatio_1 = \frac{12 \times 49}{53 \times 16} = 0.6934,$$

$$LogOddsRatio_1 = \ln(0.6934) = -0.3662,$$

and

$$V_{LogOddsRatio_1} = \frac{1}{12} + \frac{1}{53} + \frac{1}{16} + \frac{1}{49} = 0.1851.$$

This procedure is repeated for all six studies.

### Compute the summary effect using the fixed-effect model

The effect size and its variance (in log units) are copied into Table 14.5 where they are assigned the generic labels  $Y$  and  $V_Y$ .

For Saint

$$W_1 = \frac{1}{0.1851} = 5.4021,$$

$$W_1 Y_1 = 5.4021 \times (-0.3662) = -1.9780,$$

and so on for the other five studies.

The sum of  $W$  is 42.248 and the sum of  $WY$  is -30.594. From these numbers we can compute the summary effect and related statistics as follows (see (11.3) to (11.10)). In the computations that follow we use the generic  $M$  to represent the log odds ratio.

**Table 14.5** Dataset 2 – Part B (fixed-effect computations).

Study	Effect size $\gamma$	Variance within $V_\gamma$	Weight $W$	Calculated quantities		
				$W\gamma$	$W\gamma^2$	$W^2$
Saint	-0.366	0.185	5.402	-1.978	0.724	29.184
Kelly	-0.288	0.290	3.453	-0.993	0.286	11.925
Pilbeam	-0.384	0.156	6.427	-2.469	0.948	41.300
Lane	-1.322	0.058	17.155	-22.675	29.971	294.298
Wright	-0.417	0.282	3.551	-1.480	0.617	12.607
Day	-0.159	0.160	6.260	-0.998	0.159	39.190
Sum			42.248	-30.594	32.705	428.503

$$M = \frac{-30.594}{42.248} = -0.7241,$$

$$V_M = \frac{1}{42.248} = 0.0237,$$

$$SE_M = \sqrt{0.0237} = 0.1539,$$

$$LL_M = (-0.7241) - 1.96 \times 0.1539 = -1.0257,$$

$$UL_M = (-0.7241) + 1.96 \times 0.1539 = -0.4226,$$

and

$$Z = \frac{-0.7241}{0.1539} = -4.7068.$$

For a one-tailed test the p-value is given by

$$p = 1 - \Phi(-4.7068) < 0.0001,$$

and for a two-tailed test, by

$$p = 2[1 - \Phi(|-4.7068|)] < 0.0001.$$

We can convert the log odds ratio and confidence limits to the odds ratio scale using

$$OddsRatio = \exp(-0.7241) = 0.4847,$$

$$LL_{OddsRatio} = \exp(-1.0257) = 0.3586,$$

and

$$UL_{OddsRatio} = \exp(-0.4226) = 0.6553.$$

In words, using fixed-effect weights, the summary odd ratio is 0.48 with a 95% confidence interval of 0.36 to 0.66. The Z-value is -4.71, and the p-value is <0.0001 (one-tailed) or <0.0001 (two-tailed). These results are illustrated in Figure 14.3.

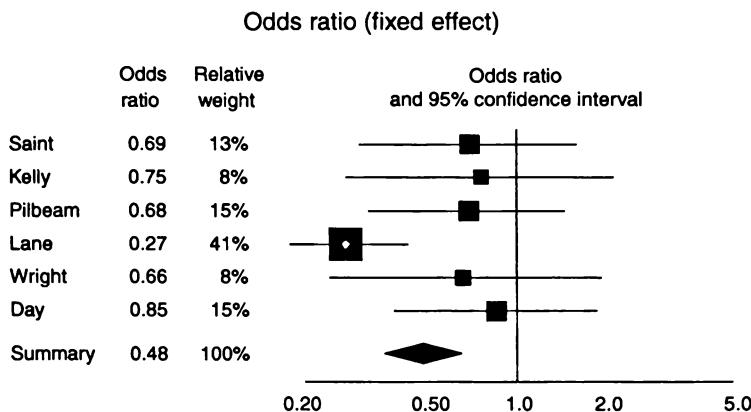


Figure 14.3 Forest plot of Dataset 2 – fixed-effect weights.

### Compute an estimate of $\tau^2$

To estimate  $\tau^2$ , the variance of the true log odds ratios, we use the DerSimonian and Laird method (see (12.2) to (12.5)). Using sums from Table 14.5,

$$Q = 32.705 - \left( \frac{-30.594^2}{42.248} \right) = 10.5512,$$

$$df = (6 - 1) = 5,$$

$$C = 42.248 - \left( \frac{428.503}{42.248} \right) = 32.1052,$$

and,

$$T^2 = \frac{10.5512 - 5}{32.1052} = 0.1729.$$

These values are reported only on a log scale.

### Compute the summary effect using the random-effects model

To compute the summary effect using the random-effects model, we use the same formulas as for the fixed effect, but the variance for each study is now the sum of the variance within studies plus the variance between studies (see (12.6) to (12.13)).

For Saint,

$$W_1^* = \frac{1}{(0.1851 + 0.1729)} = \frac{1}{(0.3580)} = 2.7932,$$

and so on for the other studies as shown in Table 14.6. Note that the within-study variance is unique for each study, but there is only one value of  $\tau^2$ , so this value (estimated as 0.173) is applied to all studies.

**Table 14.6** Dataset 2 – Part C (random-effects computations).

Study	Effect size $\gamma$	Variance within $V_\gamma$	Variance between $\tau^2$	Variance total $V_\gamma + \tau^2$	Weight $W^*$	Calculated quantities $W^* \gamma$
Saint	-0.366	0.185	0.173	0.358	2.793	-1.023
Kelly	-0.288	0.290	0.173	0.462	2.162	-0.622
Pilbeam	-0.384	0.156	0.173	0.329	3.044	-1.169
Lane	-1.322	0.058	0.173	0.231	4.325	-5.717
Wright	-0.417	0.282	0.173	0.455	2.200	-0.917
Day	-0.159	0.160	0.173	0.333	3.006	-0.479
Sum					17.531	-9.928

Then,

$$M^* = \frac{-9.928}{17.531} = -0.5663, \quad (14.3)$$

$$V_{M^*} = \frac{1}{17.531} = 0.0570, \quad (14.4)$$

$$SE_{M^*} = \sqrt{0.0570} = 0.2388,$$

$$LL_M = (-0.5663) - 1.96 \times 0.2388 = -1.0344,$$

$$UL_M = (-0.5663) + 1.96 \times 0.2388 = -0.0982,$$

$$Z^* = \frac{-0.5663}{0.2388} = -2.3711,$$

and, for a one-tailed test

$$p^* = 1 - \Phi(-2.3711) = 0.0089$$

or, for a two-tailed test

$$p^* = 2[1 - \Phi(|-2.3711|)] = 0.0177$$

We can convert the log odds ratio and confidence limits to the odds ratio scale using

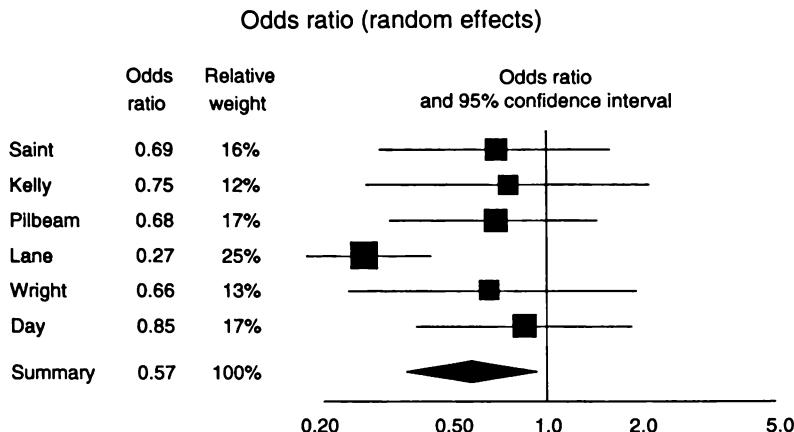
$$OddsRatio^* = \exp(-0.5663) = 0.5676,$$

$$LL_{OddsRatio^*} = \exp(-1.0344) = 0.3554,$$

and

$$UL_{OddsRatio^*} = \exp(-0.0982) = 0.9065.$$

In words, using random-effects weights, the summary odds ratio is 0.57 with a 95% confidence interval of 0.36 to 0.91. The Z-value is -2.37, and the p-value is 0.0089 (one-tailed) or 0.0177 (two-tailed). These results are illustrated in Figure 14.4.



**Figure 14.4** Forest plot of Dataset 2 – random-effects weights.

## WORKED EXAMPLE FOR CORRELATIONAL DATA (PART 1)

### Summary data

In this example we start with the correlation and sample size in six studies, as shown in Table 14.7.

**Table 14.7** Dataset 3 – Part A (basic data).

Study	Correlation	N
Fonda	0.50	40
Newman	0.60	90
Grant	0.40	25
Granger	0.20	400
Milland	0.70	60
Finch	0.45	50

### Compute the effect size and its variance for each study

For correlations, all computations are carried out using the Fisher's  $z$  transformed values (see formulas (6.2) to (6.3)). For the first study (Fonda) we compute the Fisher's  $z$  value and its variance as

$$z_1 = 0.5 \times \ln \left( \frac{1 + 0.50}{1 - 0.50} \right) = 0.5493,$$

and

$$V_1 = \frac{1}{40 - 3} = 0.0270.$$

This procedure is repeated for all six studies.

### Compute the summary effect using the fixed-effect model

The effect size and its variance (in the Fisher's  $z$  metric) are copied into Table 14.8 where they are assigned the generic labels  $Y$  and  $V_Y$ .

For Fonda

$$W_1 = \frac{1}{0.0270} = 37.0000,$$

$$W_1 Y_1 = 37.000 \times (0.5493) = 20.3243,$$

and so on for the other five studies.

The sum of  $W$  is 647.000 and the sum of  $WY$  is 242.650. From these numbers we can compute the summary effect and related statistics as follows (see (11.3) to (11.10)). In the computations that follow we use the generic  $M$  to represent the Fisher's  $z$  score.

$$M = \frac{242.650}{647.000} = 0.3750,$$

$$V_M = \frac{1}{647.000} = 0.0015,$$

$$SE_M = \sqrt{0.0015} = 0.0393,$$

$$LL_M = 0.3750 - 1.96 \times 0.0393 = 0.2980,$$

$$UL_M = 0.3750 + 1.96 \times 0.0393 = 0.4521,$$

and

$$Z = \frac{0.3750}{0.0393} = 9.5396.$$

For a one-tailed test the  $p$ -value is given by

$$p = 1 - \Phi(9.5396) < 0.0001,$$

and for a two-tailed test, by

$$p = 2[1 - \Phi(|9.5396|)] < 0.0001.$$

**Table 14.8** Dataset 3 – Part B (fixed-effect computations).

Study	Effect size $Y$	Variance within $V_Y$	Weight $W$	Calculated quantities		
				$WY$	$WY^2$	$W^2$
Fonda	0.5493	0.0270	37.000	20.324	11.164	1369.000
Newman	0.6931	0.0115	87.000	60.304	41.799	7569.000
Grant	0.4236	0.0455	22.000	9.320	3.949	484.000
Granger	0.2027	0.0025	397.000	80.485	16.317	157609.000
Milland	0.8673	0.0175	57.000	49.436	42.876	3249.000
Finch	0.4847	0.0213	47.000	22.781	11.042	2209.000
Sum			647.000	242.650	127.147	172489.000

We can convert the effect size and confidence limits from the Fisher's  $z$  metric to correlations using

$$r = \frac{e^{(2 \times 0.3750)} - 1}{e^{(2 \times 0.3750)} + 1} = 0.3584,$$

$$LL_r = \frac{e^{(2 \times 0.2980)} - 1}{e^{(2 \times 0.2980)} + 1} = 0.2895,$$

and

$$UL_r = \frac{e^{(2 \times 0.4521)} - 1}{e^{(2 \times 0.4521)} + 1} = 0.4236.$$

In words, using fixed-effect weights, the summary estimate of the correlation is 0.36 with a 95% confidence interval of 0.29 to 0.42. The  $Z$ -value is 9.54, and the  $p$ -value is <0.0001 (one-tailed) or <0.0001 (two tailed). These results are illustrated in Figure 14.5.

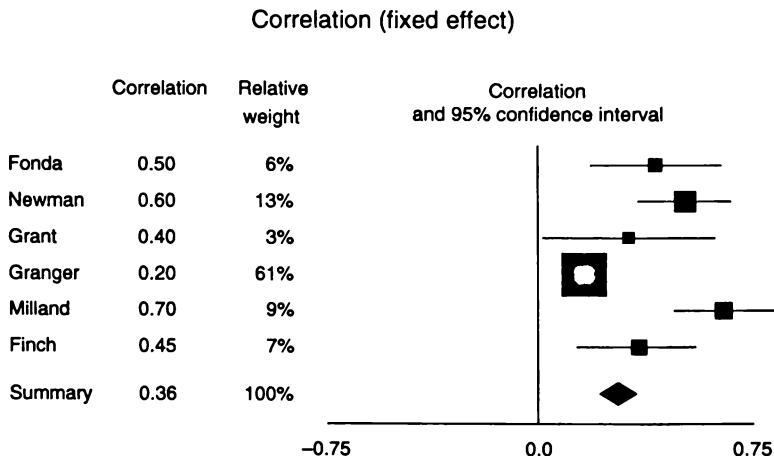


Figure 14.5 Forest plot of Dataset 3 – fixed-effect weights.

### Compute an estimate of $\tau^2$

To estimate  $\tau^2$ , the variance of the true Fisher's  $z$ , we use the DerSimonian and Laird method (see (12.2) to (12.5)). Using sums from Table 14.8,

$$Q = 127.147 - \left( \frac{242.650^2}{647.000} \right) = 36.1437,$$

$$df = (6 - 1) = 5,$$

$$C = 647.000 - \left( \frac{172489.000}{647.000} \right) = 380.4019,$$

and

$$T^2 = \frac{36.1437 - 5}{380.4019} = 0.0819.$$

### Compute the summary effect using the random-effects model

To compute the summary effect using the random-effects model, we use the same formulas as for the fixed effect, but the variance for each study is now the sum of the variance within studies plus the variance between studies (see (12.6) to (12.13)).

For Fonda,

$$W_i^* = \frac{1}{(0.0270 + 0.0819)} = \frac{1}{(0.1089)} = 9.1829,$$

and so on for the other studies as shown in Table 14.9. Note that the within-study variance is unique for each study, but there is only one value of  $\tau^2$ , so this value (estimated as 0.0819) is applied to all studies.

Then,

$$M^* = \frac{31.621}{59.351} = 0.5328, \quad (14.5)$$

$$V_{M^*} = \frac{1}{59.351} = 0.0168, \quad (14.6)$$

$$SE_{M^*} = \sqrt{0.0168} = 0.1298,$$

$$LL_{M^*} = (0.5328) - 1.96 \times 0.1298 = 0.2784,$$

$$UL_{M^*} = (0.5328) + 1.96 \times 0.1298 = 0.7872,$$

and

$$Z^* = \frac{0.5328}{0.1298} = 4.1045.$$

Then, for a one-tailed test

$$p^* = 1 - \Phi(4.1045) < 0.0001,$$

or, for a two-tailed test

$$p^* = 2[1 - \Phi(|4.1045|)] < 0.0001.$$

**Table 14.9** Dataset 3 – Part C (random-effects computations).

Study	Effect size $Y$	Variance within $V_Y$	Variance between $T^2$	Variance total $V_Y + T^2$	Weight $W^*$	Calculated quantities $W^* Y$
Fonda	0.549	0.027	0.082	0.109	9.183	5.044
Newman	0.693	0.012	0.082	0.093	10.711	7.424
Grant	0.424	0.046	0.082	0.127	7.854	3.327
Granger	0.203	0.003	0.082	0.084	11.850	2.402
Milland	0.867	0.018	0.082	0.099	10.059	8.724
Finch	0.485	0.021	0.082	0.103	9.695	4.699
Sum					59.351	31.621

We can convert the effect size and confidence limits from the Fisher's  $z$  metric to correlations using

$$r^* = \frac{e^{(2 \times 0.5328)} - 1}{e^{(2 \times 0.5328)} + 1} = 0.4875,$$

$$LL_{r^*} = \frac{e^{(2 \times 0.2784)} - 1}{e^{(2 \times 0.2784)} + 1} = 0.2714,$$

and

$$UL_{r^*} = \frac{e^{(2 \times 0.7872)} - 1}{e^{(2 \times 0.7872)} + 1} = 0.6568.$$

In words, using random-effects weights, the summary estimate of the correlation is 0.49 with a 95% confidence interval of 0.27 to 0.66. The Z-value is 4.10, and the p-value is <0.0001 (one-tailed) or <0.0001 (two tailed). These results are illustrated in Figure 14.6.

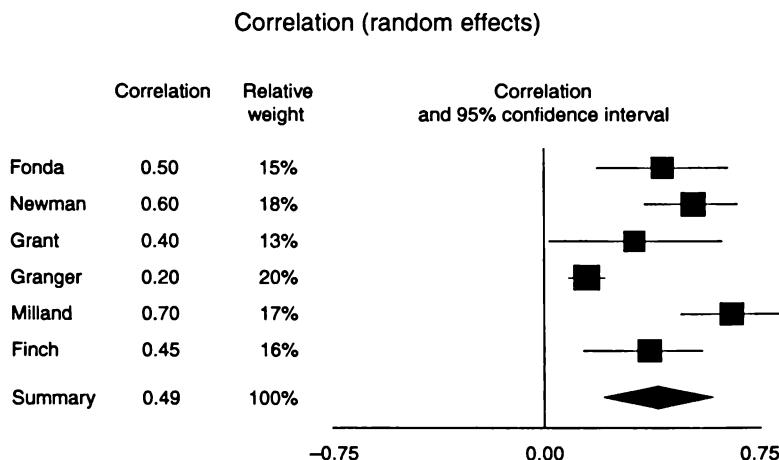


Figure 14.6 Forest plot of Dataset 3 – random-effects weights.

### SUMMARY POINTS

- This chapter includes worked examples showing how to compute the summary effect using fixed-effect and random-effects models.
- For the standardized mean difference we work with the effect sizes directly.
- For ratios we work with the log transformed data.
- For correlations we work with the Fisher's  $z$  transformed data.
- These worked examples are available as Excel files on the book's website ([www.Introduction-to-Meta-Analysis.com](http://www.Introduction-to-Meta-Analysis.com)).

# Heterogeneity



# Overview

- 
- Introduction
  - Nomenclature
  - Worked examples
- 

## INTRODUCTION

A central theme in this volume is that the goal of a synthesis is not simply to compute a summary effect, but rather to make sense of the pattern of effects. An intervention that consistently reduces the risk of criminal behavior by 40% across a range of studies is very different from one that reduces the risk by 40% on average with a risk reduction that ranges from 10% in some studies to 70% in others. If the effect size is consistent across studies we need to know that and to consider the implications, and if it varies we need to know that and consider the different implications. This issue, the heterogeneity in effect sizes, is the focus of the next two Parts.

In this Part we show how to identify and quantify the heterogeneity in effect sizes. There are a number of questions we might ask about heterogeneity, and there are statistics that address each question. Each statistic addresses a specific question, and these statistics are not interchangeable. It is important to understand what each statistic tells us (and does not tell us). The following are some of the common questions, and the corresponding statistic

- Is there evidence of heterogeneity in true effects sizes? This is addressed by the  $Q$  value, the degrees of freedom, and the corresponding p-value.
- Of the variance that we see (in the observed effects) what proportion reflects variance in true effects rather than sampling error. This is addressed by the  $I^2$  statistic
- How much do the true effects vary? This is addressed by  $T^2$  (the variance of true effects) and  $T$  (the standard deviation of true effects). Note that these statistics may be presented on the log scale (or in another metric)
- What is the interval of true effects? This is addressed by the prediction interval. This takes into account information about the mean as well as the heterogeneity. Critically, the interval is reported on the same scale as the effect size.

Researchers often conflate these statistics with each other, and use them in ways for which they were not intended. To address this we include a chapter that outlines how we should think about heterogeneity. Once we are clear about what we intend to ask, it becomes clear what statistics actually address that question.

Finally, we include a chapter to address a serious mistake that is common in the literature. Researchers often classify heterogeneity as being low, moderate or high based on the value of  $I^2$ . We explain why this approach should be avoided.

Where Part 4 focuses on ways to quantify heterogeneity, the next part (Part 5) introduces methods to explore the reasons for heterogeneity. In Chapter 21 we show how to compare the effect size in different subgroups of studies (such as studies that used different populations, or studies that used different variants of the intervention), similar to analysis of variance in primary studies. In Chapter 22 we show how to assess the relationship between effect size and covariates (such as the dose of a drug or the mean age in the study sample), similar to multiple regression in primary studies. Finally, Chapter 23 includes a discussion of issues (and some important caveats) related to both techniques.

## NOMENCLATURE

As we did in the previous part we distinguish between a true effect size and an observed effect size. A study's *true effect size* is the effect size in the underlying population, and is the effect size that we would observe if the study had an infinitely large sample size (and therefore no sampling error). A study's *observed effect size* is the effect size that is actually observed.

We use the terms *variation* and *dispersion* to refer to differences among values, sometimes true effects and sometimes observed effects, depending on the context. By contrast, we use *heterogeneity* to mean heterogeneity in true effects only.

We shall introduce several measures of heterogeneity, one of which is tau-squared, defined as the variance of the true effect sizes. We will use  $\tau^2$  to refer to the parameter (the population value) and  $T^2$  to refer to our estimate of this parameter.

## WORKED EXAMPLES

The discussion of formulas and concepts in Chapter 16 is followed by a section with worked examples using a series of effect sizes (standardized mean differences, odds ratios, and correlations). While reading the sections on formulas and concepts, it will be helpful to refer to one or more of the worked examples.

These examples are continuations of the three worked examples in Chapter 14, and are all available as Excel spreadsheets on the book's website.

# Identifying and Quantifying Heterogeneity

---

Introduction

Isolating the variation in true effects

Computing  $Q$

Estimating  $\tau^2$

The  $\ell^2$  statistic

Comparing the measures of heterogeneity

Confidence intervals for  $\tau^2$

Confidence intervals (or uncertainty intervals) for  $\ell^2$

---

## INTRODUCTION

Under the random-effects model we allow that the true effect size may vary from study to study. In this chapter we discuss approaches to identify and then quantify this heterogeneity.

## ISOLATING THE VARIATION IN TRUE EFFECTS

The mechanisms for describing the variation among scores in a primary study are well known. We can compute the standard deviation of the scores and discuss the proportion of subjects falling within a given range. We can compute the variance of the scores and discuss what proportion of variance can be explained by covariates.

Our goals are similar in a meta-analysis, in the sense that we want to describe the variation, using indices such as the standard deviation and variance. However, the process is more complicated for the following reason. When we speak about the *heterogeneity* in effect sizes we mean the variation in the *true effect sizes*. However, the variation that we actually observe is partly spurious, incorporating both (true) heterogeneity and also random error.

To understand the problem, suppose for a moment that all studies in the analysis shared the same *true effect size*, so that the (true) heterogeneity is zero. Under this assumption, we would not expect the observed effects to be identical to each other. Rather, because of within-study error, we would expect each to fall *within some range* of the common effect.

Now, assume that the true effect size *does* vary from one study to the next. In this case, the observed effects vary from one another for two reasons. One is the real heterogeneity in effect size, and the other is the within-study error. If we want to quantify the heterogeneity we need to partition the observed variation into these two components, and then focus on the former.

The mechanism that we use to extract the true between-studies variation from the observed variation is as follows.

1. We compute the total amount of study-to-study variation actually observed.
2. We estimate how much the observed effects would be expected to vary from each other if the true effect was actually the same in all studies.
3. The excess variation (if any) is assumed to reflect real differences in effect size (that is, the heterogeneity).

Consider the top row in Figure 16.1. The observed effects (and therefore variation in the observed effects) are identical in A and B. The difference between A and B is that the confidence intervals for each study in A are relatively wide, while the confidence intervals for each study in B are relatively narrow. The visual impression in A is that all studies could share a common effect, with the observed dispersion falling within

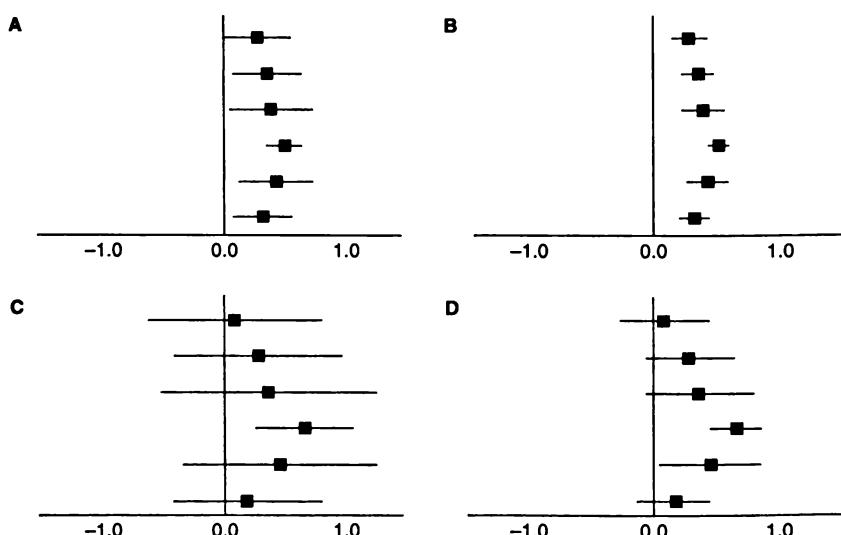


Figure 16.1 Dispersion across studies relative to error within studies.

the umbrella of the confidence intervals. By contrast, the confidence intervals for the  $B$  studies are quite narrow, and cannot comfortably account for the observed dispersion.

Similarly, consider the bottom row in this plot. Again, the *observed* effects are identical in  $C$  and  $D$ , but  $C$  has wider confidence intervals around each study. In  $C$ , the effects can be fully explained by within-study error, while in  $D$  they cannot.

The difference between  $A$  and  $B$  (on the one hand) versus  $C$  and  $D$  (on the other) is that we have changed the *absolute value* (or the scale) of the dispersion. To move from  $A$  to  $C$  we multiplied both the within-study variance and the observed variance by 2.0 so that the scale has increased but the ratio (observed/within) is unchanged. The same holds true for  $B$  and  $D$ . While the effects are more widely dispersed in the second row than in the first, this is not relevant to the purpose of isolating the true dispersion. What matters is only the *ratio* of observed to expected dispersion, which is the same in  $A$  and  $C$  (and is the same in  $B$  and  $D$ ).

It is this aspect of the dispersion – the one reflected in the difference of  $A$  versus  $B$ , and of  $C$  versus  $D$  that we want to capture for the purpose of isolating the true variation. In other words, we need a statistic that is sensitive to the ratio of the observed variation to the within-study error, rather than their absolute values. The statistic that we use for this purpose is  $Q$ .

## COMPUTING Q

The first step in partitioning the variation is to compute  $Q$ , defined as

$$Q = \sum_{i=1}^k W_i(Y_i - M)^2, \quad (16.1)$$

where  $W_i$  is the study weight ( $1/V_i$ ),  $Y_i$  is the study effect size, and  $M$  is the summary effect and  $k$  is the number of studies. In words, we compute the deviation of each effect size from the mean, square it, weight this by the inverse-variance for that study, and sum these values over all studies to yield the weighted sum of squares (WSS), or  $Q$ .

The same formula can be written as

$$Q = \sum_{i=1}^k \left( \frac{Y_i - M}{S_i} \right)^2 \quad (16.2)$$

to highlight the fact that  $Q$  is a standardized measure, which means that it is not affected by the metric of the effect size index. The analogy would be to the standardized mean difference  $d$ , where the mean difference is divided by the within-study standard deviation.

Finally, an equivalent formula, useful for computations, is

$$Q = \sum_{i=1}^k W_i Y_i^2 - \frac{\left( \sum_{i=1}^k W_i Y_i \right)^2}{\sum_{i=1}^k W_i}. \quad (16.3)$$

### The expected value of $Q$ based on within-study error

The next step is to determine the expected value of  $Q$  on the assumption that all studies share a common effect size, and (it follows) all the variation is due to sampling error within studies. Because  $Q$  is a standardized measure the expected value does not depend on the metric of effect size, but is simply the degrees of freedom ( $df$ ),

$$df = k - 1, \quad (16.4)$$

where  $k$  is the number of studies.

### The excess variation

Since  $Q$  is the observed WSS and  $df$  is the expected WSS (under the assumption that all studies share a common effect), the difference,

$$Q - df,$$

reflects the excess variation, the part that will be attributed to differences in the true effects from study to study.

### Ratio of observed to expected variation

Earlier, we used Figure 16.1 to introduce the concept of excess variation, and showed that it needs to be based on the ratio of observed variance to the within-study variance. In Figure 16.2 we reproduce the same four plots but add the  $Q$  statistic, which quantifies this concept. The plot also includes additional statistics, which will be explained below.

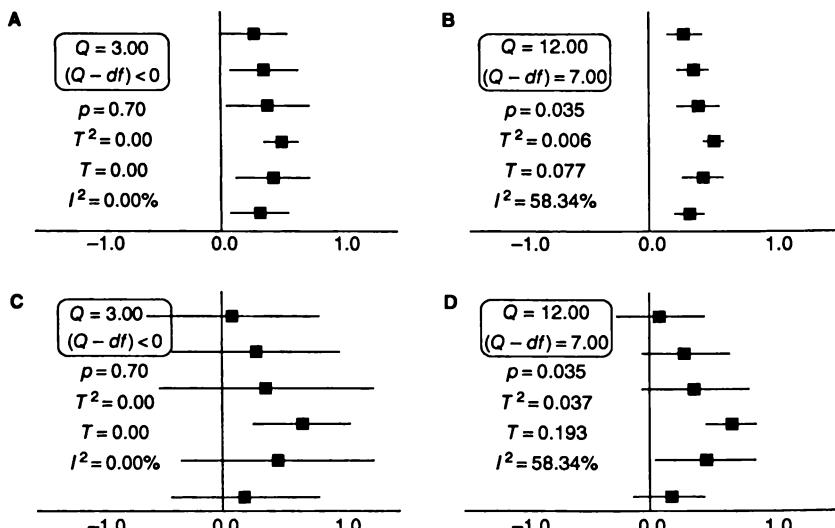


Figure 16.2  $Q$  in relation to  $df$  as measure of dispersion.

First, consider the top row. For plot A the observed value of  $Q$  is 3.00, versus an expected value (under the assumption of a common effect size) of 5.00 (that is,  $k - 1$ ). In this case the observed variation is less than we would expect based on within-study error ( $Q$  is less than the degrees of freedom). For plot B the observed value of  $Q$  is 12.00, versus an expected value of 5.00, so the observed variation is greater than we would expect based on within-study error ( $Q$  is greater than the degrees of freedom). This, of course, is consistent with the visual impression discussed earlier. We see the same thing in the second row, where for plot C the observed value of  $Q$  is 3.00, versus an expected value of 5.00, ( $Q < df$ ) and for plot D the observed value of  $Q$  is 12.00, versus an expected value of 5.00, ( $Q > df$ ).

Note that  $Q$  is the same (3.00) in A and C because these plots share the same ratio, despite the fact that the absolute range of effects is higher in C. Similarly, the value of  $Q$  is the same (12.00) in B and D because these plots share the same ratio, despite the fact that the absolute range of effects is higher in D.

At this point we have the values  $Q$ , which reflects the total dispersion (WSS) and  $Q - df$ , which reflects the excess dispersion. However,  $Q$  is not an intuitive measure. For one thing,  $Q$  is a sum (rather than a mean) and, as such, depends strongly on the number of studies. For another,  $Q$  is on a standardized scale, and for some purposes we will want to express the dispersion either as a ratio or on the same scale as the effect size itself.

Now that we have  $Q$ , however, we can use it to construct measures that do address specific needs, as outlined in Figure 16.3. To test the assumption of homogeneity we will work directly with  $Q$ , and take advantage of the fact that it is on a standardized scale and is sensitive to the number of studies. To estimate the variance (and standard deviation) of the true effects we will start with  $Q$ , remove the dependence on the number of studies, and return to the original metric. These estimates are called  $T^2$  and  $T$ . Finally, to estimate what proportion of the observed variance reflects real differences among studies (rather than random error) we will start with  $Q$ , remove the dependence on the number of studies, and express the result as a ratio (called  $I^2$ ).

### **Testing the assumption of homogeneity in effects**

Researchers typically ask if the heterogeneity is statistically significant, and we can use  $Q$  (and  $df$ ) to address this question. Formally, we pose the null hypothesis that all studies share a common effect size and then test this hypothesis. Under the null hypothesis,  $Q$  will follow a central chi-squared distribution with degrees of freedom equal to  $k - 1$ , so we can report a  $p$ -value for any observed value of  $Q$ . Typically, we set alpha at 0.10 or at 0.05, with a  $p$ -value less than alpha leading us to reject the null hypothesis, and conclude that the studies do not share a common effect size.

This test of significance, like all tests of significance, is sensitive both to the magnitude of the effect (here, the excess dispersion) and the precision with which this effect is estimated (here, based on the number of studies).

The impact of the excess dispersion on the  $p$ -value is evident if we compare plots A versus B in Figure 16.4. These plots both have six studies, and as the excess dispersion increases ( $Q$  moves from 3.00 in A to 12.00 in B) the  $p$ -value moves from 0.70 to 0.035.

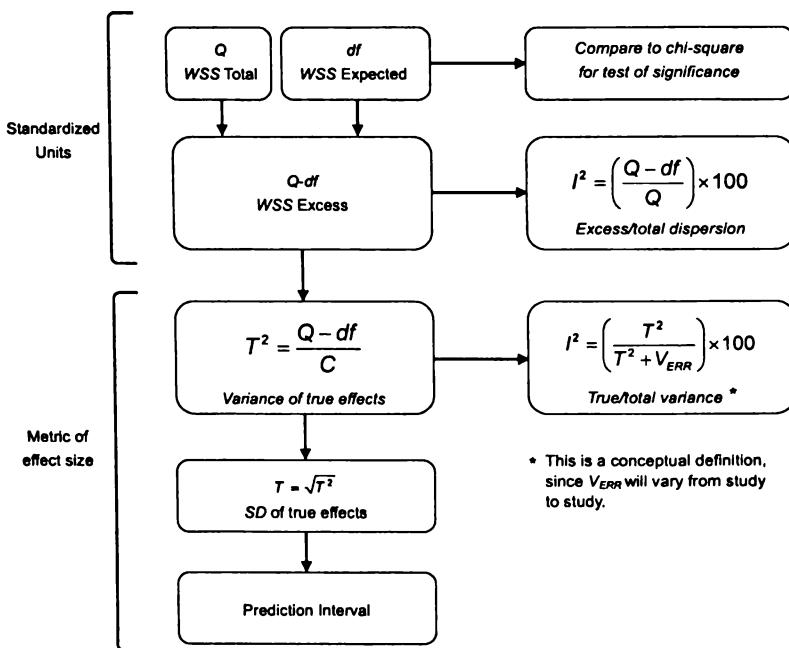


Figure 16.3 Flowchart showing how  $T^2$  and  $I^2$  are derived from  $Q$  and  $df$ .

Similarly, compare plots C and D, which both have twelve studies. As we move from C to D,  $Q$  moves from 6.62 to 27.12 and the  $p$ -value moves from 0.830 to 0.004.

The impact of the number of studies is evident if we compare plots A versus C. The two plots are essentially identical, except that A has six studies and C has twelve studies with the same estimated value of between-studies variation (prior to being truncated at zero). With the additional precision the  $p$ -value moves away from zero, from 0.70 (for A) to 0.83 (for C).

Similarly, the impact of the number of studies is evident if we compare plots B versus D. The two plots are essentially identical, except that (again) B has six studies and D has twelve studies with the same estimated value of between-studies variation ( $T^2 = 0.037$ ). With the additional precision the  $p$ -value moves towards zero, from 0.035 (for B) to 0.004 (for D).

Note that the  $p$ -value for the left-hand columns moved toward 1.0 as we added studies, while the  $p$ -value for the right-hand columns moved toward 0.0 as we added studies. At left, since  $Q$  is less than  $df$  the additional evidence strengthens the case that the excess dispersion is zero, and moves the  $p$ -value towards 1.0. At right, since  $Q$  exceeds  $df$ , the additional evidence strengthens the case that the excess dispersion is *not* zero, and moves the  $p$ -value towards 0.0.

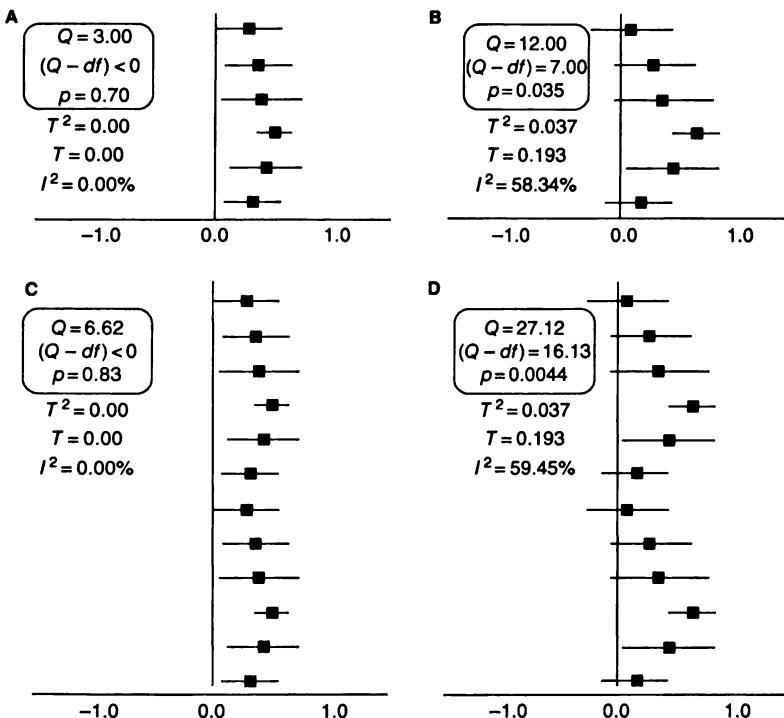


Figure 16.4 Impact of  $Q$  and number of studies on the  $p$ -value.

### Some remarks about $Q$ and its $p$ -value

The test of the null hypothesis (that all studies share a common effect size) is subject to the same caveats as all tests of significance, as follows.

First, while a significant  $p$ -value provides evidence that the true effects vary, the converse is not true. A nonsignificant  $p$ -value should not be taken as evidence that the effect sizes are consistent, since the lack of significance may be due to low power. With a small number of studies and/or large within-study variance (small studies), even substantial between-studies dispersion might yield a nonsignificant  $p$ -value.

Second, the  $Q$  statistic and  $p$ -value address only the test of significance and should never be used as surrogates for the amount of true variance. A nonsignificant  $p$ -value could reflect a trivial amount of observed dispersion, but could also reflect a substantial amount of observed dispersion with imprecise studies. Similarly, a significant  $p$ -value could reflect a substantial amount of observed dispersion, but could also reflect a minor amount of observed dispersion with precise studies.

In sum, the purpose served by this test is to assess the viability of the null hypothesis, and not to estimate the magnitude of the true dispersion. There are several ways that we can describe the dispersion of true effect sizes, and we shall deal with each of these in the balance of this chapter.

## ESTIMATING $\tau^2$

The parameter tau-squared ( $\tau^2$ ) is defined as the variance of the true effect sizes. In other words, if we had an infinitely large sample of studies, each, itself, infinitely large (so that the estimate in each study was the true effect) and computed the variance of these effects, this variance would be  $\tau^2$ .

Since we cannot observe the true effects we cannot compute this variance directly. Rather, we estimate it from the observed effects, with the estimate denoted  $T^2$ . To yield this estimate we start with the difference ( $Q - df$ ) which represents the dispersion in true effects on a standardized scale. We divide by a quantity ( $C$ ) which has the effect of putting the measure back into its original metric and also of making it an average, rather than a sum, of squared deviations. Concretely,

$$T^2 = \frac{Q - df}{C} \quad (16.5)$$

where

$$C = \sum W_i - \frac{\sum W_i^2}{\sum W_i}. \quad (16.6)$$

This means that  $T^2$  is in the same metric (squared) as the effect size itself, and also reflects the absolute amount of variation in that scale.

While the actual variance of the true effects ( $\tau^2$ ) can never be less than zero, our estimate of this value ( $T^2$ ) can be less than zero if, because of sampling error, the observed variance is less than we would expect based on within-study error – in other words, if  $Q < df$ . In this case,  $T^2$  is simply set to zero.

If  $Q > df$  then  $T^2$  will be positive, and it will be based on two factors. The first is the amount of excess variation ( $Q - df$ ), and the second is the metric of the effect size index.

The impact of the excess variation on our estimate of  $T^2$  is evident if we compare plots A versus B in Figure 16.5. The within-study error is smaller in B. Therefore, while the observed variation is the same in both plots, a higher proportion of this variation is assumed to be real in B. As we move from A to B,  $Q$  moves from 12.00 to 48.01, and  $T^2$  from 0.037 to 0.057.

The impact of the scale on our estimate of  $T^2$  is evident if we compare plots C versus D in Figure 16.5.  $Q$  and  $df$  are the same in the two plots, which means that the same proportion of the observed variance will be attributed to between-studies variance. However, the absolute amount of the variance is larger in D, so this proportion translates into a larger estimate of  $\tau^2$ . As we move from C to D,  $T^2$  moves from 0.037 to 0.096.

### **Some remarks about $T^2$**

As our estimate for the variance of the true effects,  $T^2$  is used to assign weights under the random-effects model, where the weight assigned to each study is

$$W_i^* = \frac{1}{V_{Y_i}^*} = \frac{1}{V_{Y_i} + T^2}. \quad (16.7)$$

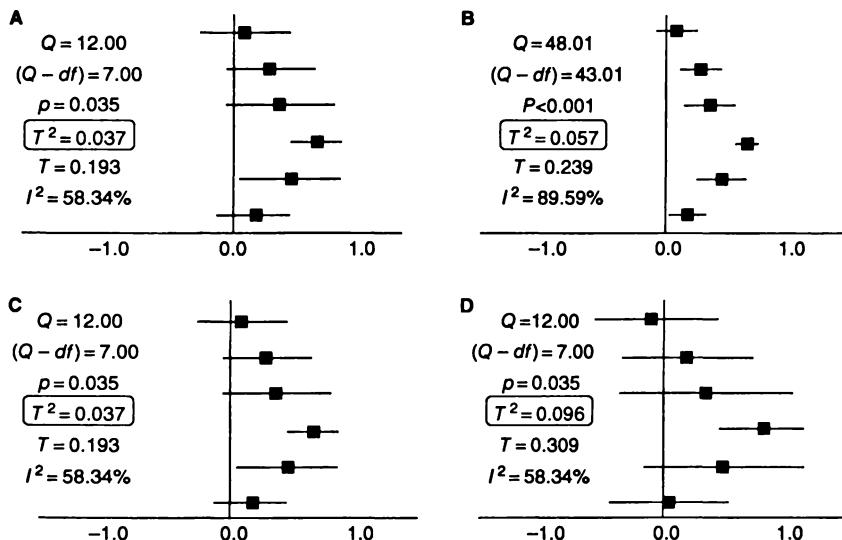


Figure 16.5 Impact of excess dispersion and absolute dispersion on  $T^2$ .

In words, the total variance for a study ( $V_{Y^*}$ ) is the sum of the within-study variance ( $V_Y$ ) and the between-studies variance ( $T^2$ ).

This method of estimating the variance between studies is the most popular, and is known as the *method of moments* or the *DerSimonian and Laird method*. This method does not make any assumptions about the distribution of the random effects. It also has the advantage of being the easiest to compute and the easiest to explain, which makes it useful for a text. Alternatives exist, and some statisticians favor a restricted maximum likelihood (REML) method.

Interestingly, one of the authors of the key DerSimonian and Laird paper has since argued that the simple method we describe should no longer be used, since computational simplicity is no longer an important consideration. However, we believe that it is instructive to describe this simple method, and note that differences in results from one method to the other are likely to be small. Formulas to compute confidence intervals for  $T^2$  are presented at the end of this chapter.

### Tau

Above, we discussed the variance of the true effect sizes, where  $\tau^2$  refers to the actual variance and  $T^2$  is our estimate of this parameter. Now, we turn to the standard deviation of the true effect sizes. Here,  $\tau$  refers to the actual standard deviation and  $T$  is our estimate of this parameter.

$T$ , the estimate of the standard deviation, is simply the square root of  $T^2$ ,

$$T = \sqrt{T^2}. \quad (16.8)$$

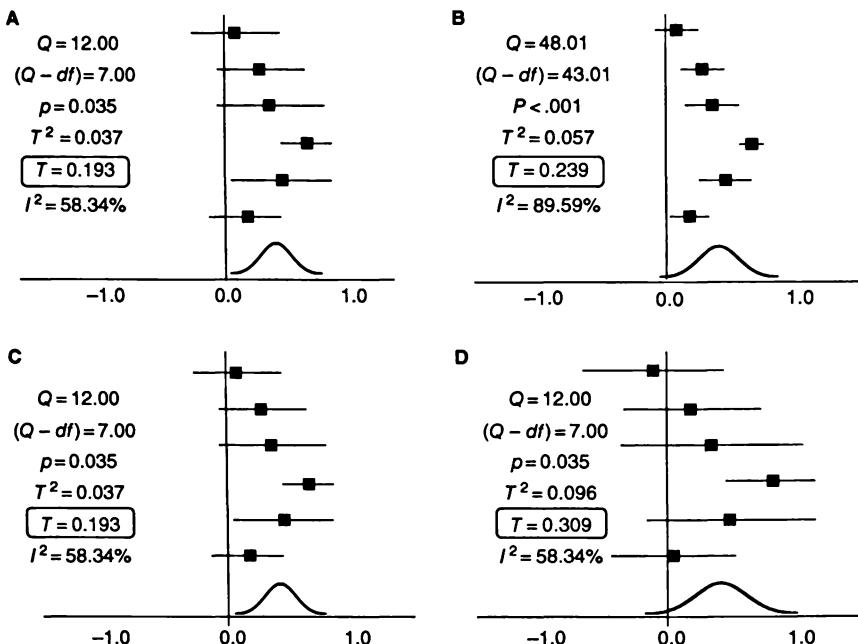


Figure 16.6 Impact of excess and absolute dispersion on  $T$ .

Like  $T^2$ ,  $T$  is on the same scale as the effect size itself, but  $T^2$  (a variance) is a squared value, while  $T$  (a standard deviation) is not. Like the standard deviation in a primary study,  $T$  can be used to describe the distribution of effect sizes about the mean effect. If we are willing to assume that the effects are normally distributed (and we have a reasonably precise estimate of  $T$ ), we can get a sense for the range of true effect sizes, and then consider the substantive implications of this range.

Figure 16.6 is identical to Figure 16.5 but this time we have added to each plot the expected distribution of true effects, based on  $T$ . For example, in plot A the summary effect is 0.41 and  $T$  is 0.193. We expect that some 95% of the true effects will fall in the range of 0.41 plus or minus 1.96  $T$ , or 0.04 to 0.79, and this is the range reflected in the bell curve. The same approach is used to construct the curve for all the plots.

Recall that plots A and B have the same observed variance, but differ in the proportion of this variance that is attributed to real differences in effect size. In A, the bell curve is relatively narrow, and captures only a fraction of the observed dispersion – the rest is assumed to reflect error. In B, the bell curve is relatively wide, and captures a larger fraction of the dispersion, since most of the dispersion is here assumed to be real.

Similarly, recall that in plots C and D the ratio of true to observed variance is the same, but the observed dispersion (the scale) is larger in D. The bell curve is wider in D than in C (because of the different scale), but in both cases a comparable proportion of the effects fall within the range of the curve (because the ratio is the same).

### Some remarks about $T$

Our estimate  $T$  of the standard deviation of the true effects enables us to talk about the substantive importance of the dispersion. Suppose an intervention has a summary effect size of 0.50. If  $T$  is 0.10, then most of the effects (95 %) fall in the approximate range of 0.30 to 0.70. If  $T$  is 0.20 then most of the true effects fall in the approximate range of 0.10 to 0.90. If  $T$  is 0.30 then most of the true effects fall in the approximate range of -0.10 to +1.10. This interval is called a prediction interval. We still need to attach a value judgment to these ranges (what effect size is *harmful*, what effect size is *trivial*, what effect size is *useful*), but by having a sense of the distribution we have at least a starting point for these discussions.

In this example we assume that the effect size and  $T$  are estimated accurately. In practice, if we wanted to make predictions about the distribution of true effects we would need to take account of the error in estimating both of these values. Toward the end of this chapter we show how to compute confidence intervals for the value of  $T$ , and in Chapter 17 we show how to compute prediction intervals that take account of these sources of error.

### THE $I^2$ STATISTIC

The utility of  $T^2$  and  $T$  lies in the fact that they are absolute measures, which means that they quantify deviation on the same scale as the effect size index. In some cases, however, we would prefer to think about heterogeneity independent of scale and ask *What proportion of the observed variance reflects real differences in effect size?*

Higgins *et al.* (2003) proposed using a statistic,  $I^2$ , to reflect this proportion, that could serve as a kind of signal-to-noise ratio. It is computed as

$$I^2 = \left( \frac{Q - df}{Q} \right) \times 100\%, \quad (16.9)$$

that is, the ratio of excess dispersion to total dispersion. The statistic  $I^2$  can be viewed as a statistic of the form

$$I^2 = \left( \frac{\text{Variance}_{\text{het}}}{\text{Variance}_{\text{total}}} \right) \times 100\% = \left( \frac{\tau^2}{\tau^2 + V_Y} \right) \times 100\%, \quad (16.10)$$

that is, the ratio of true heterogeneity to total variance across the observed effect estimates. However, this is not a true definition of  $I^2$  because in reality there is not a single  $V_Y$ , since the within-study variances vary from study to study. The  $I^2$  statistic is a descriptive statistic and not an estimate of any underlying quantity.

To get a sense of the factors that control (and do not control)  $I^2$ , consider Figure 16.7. For any given  $df$ ,  $I^2$  moves in tandem with  $Q$ . As such, it is driven entirely by the ratio of observed dispersion to within-study dispersion. In the top row, both plots A and B have a  $Q$  value of 12.00 with 5 degrees of freedom. Therefore, they both have the same value of  $I^2$ , 58.34%. The fact that A has a wider scale than B (which yields a higher value of  $T^2$ ) does not impact on  $I^2$ . Similarly, in the bottom row both plots C and D have a  $Q$  value of 48.01 with 5 degrees of freedom, and therefore both have the same  $I^2$ , 89.59%. Again, the fact that C has a wider scale than D does not impact on  $I^2$ .

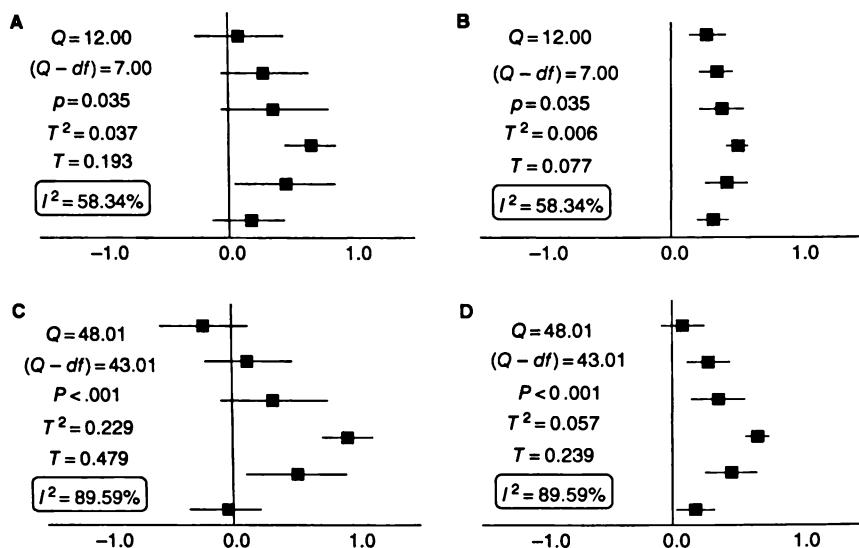


Figure 16.7 Impact of excess dispersion on  $I^2$ .

$I^2$  reflects the extent of overlap of confidence intervals, which is not dependent on the actual location or spread of the true effects. As such it is convenient to view  $I^2$  as a measure of *inconsistency* across the findings of the studies, and not as a measure of the real variation across the underlying true effects.

The scale of  $I^2$  has a range of 0–100%, regardless of the scale used for the meta-analysis itself. It can be interpreted as a ratio, and has the additional advantage of being analogous to indices used in psychometrics (where reliability is the ratio of true to total variance) or regression (where  $R^2$  is the proportion of the total variance that can be explained by the covariates). Importantly,  $I^2$  is not directly affected by the number of studies in the analysis.

#### Some remarks about $I^2$

Higgins *et al.* (2003) provide some tentative benchmarks for  $I^2$ . They suggest that values on the order of 25%, 50%, and 75% might be considered as low, moderate, and high, respectively. Some context for the interpretation of  $I^2$  is provided by a survey of meta-analyses of clinical trials in the Cochrane Database of Systematic Reviews, reported by Higgins *et al.* (2003). The value of  $I^2$  was zero for about half of the meta-analyses, and was distributed evenly between 0% to 100% for the other half. It is likely that  $I^2$  would be distributed differently in meta-analyses of other fields or other kinds of studies.

Note that the benchmarks (like the index itself) refer to the question of what *proportion* of the observed variation is real, and not to the variation on an absolute scale. An

$I^2$  value near 100% means only that most of the observed variance is real, but does not imply that the effects are dispersed over a wide range (they could fall in a narrow range but be estimated precisely). Nor does a low value of  $I^2$  imply that the effect are clustered in a narrow range (the observed effects could vary across a wide range, in studies with a lot of error). As such,  $I^2$  is not meant to address the substantive implications of the dispersion.

Formulas to compute confidence intervals for  $I^2$  are presented at the end of this chapter.

### COMPARING THE MEASURES OF HETEROGENEITY

We have described five ways of measuring heterogeneity,  $Q$ ,  $p$ ,  $T^2$ ,  $T$ , and  $I^2$ . Table 16.1 shows the relationship among these measures. Since all the indices are based on  $Q$  (in relation to  $df$ ), it follows that all will be low (or zero) if the total dispersion is low relative to the error within studies, and higher if the total dispersion is high, relative to the error within studies. However, the various measures of heterogeneity build on this core in different ways which makes each useful for a specific purpose.

Note that the estimates  $T^2$  and  $T$  are based on the excess dispersion, but the population values ( $\tau^2$  and of  $\tau$ ) are defined solely by the variance of the true effects.

- The  $Q$  statistic and its  $p$ -value serve as a test of significance. The qualities that make these useful for this purpose are that they are sensitive to the number of studies and they are not sensitive to the metric of the effect size index.
- Our estimate of  $\tau^2$  serves as the between-studies variance in the analysis and our estimate of  $\tau$  serves as the standard deviation of the true effects. The qualities that make these useful for this purpose are that they are sensitive to the metric of the effect size, and they are not sensitive to the number of studies.
- $I^2$  is the ratio of true heterogeneity to total variation in observed effects, a kind of signal to noise ratio. The qualities that make it useful for this purpose are that it is not sensitive to the metric of the effect size and it is not sensitive to the number of studies.

It is important to understand that  $T^2$  and  $T$  (on the one hand) and  $I^2$  (on the other) serve two entirely different functions. The statistics  $T^2$  (and  $T$ ) reflect the amount of true

Table 16.1 Factors affecting measures of dispersion.

	Range of possible values	Depends on number of studies	Depends on scale
$Q$	$0 \leq Q$		✓
$p$	$0 \leq p \leq 1$		✓
$T^2$	$0 \leq T^2$		
$T$	$0 \leq T$		✓
$I^2$	$0\% \leq I^2 \leq 100\%$		✓

heterogeneity (the variance or the standard deviation) while  $I^2$  reflects the proportion of observed dispersion that is due to this heterogeneity. In a sense, if we were to multiply the observed variance by  $I^2$ , we would get  $T^2$  (this is meant as an illustration only, since the actual computation is more complicated). As such, the two tend to move in tandem, but have very different meanings.

$I^2$  reflects only the proportion of variance that is true, and says nothing about the absolute value of this variance. In Figure 16.8, plots A and B have the same value of  $I^2$  (58.34%) but in A the true effects are clustered in a small range ( $T^2 = 0.006$ ) while in B they are dispersed across a wider range ( $T^2 = 0.037$ ).

Conversely,  $T^2$  reflects only the absolute value of the true variance and says nothing about the proportion of observed variance that is true. In Figure 16.9,  $T^2$  is the same in both plots, but in A it is a large part ( $I^2 = 58.34\%$ ) of a small observed dispersion whereas in B it is a small part ( $I^2 = 16.01\%$ ) of a large observed dispersion.

Note also that  $T^2$  is tied to the effect size index while  $I^2$  is not. For example,  $T^2$  for a synthesis of risk ratios will be in the metric of log risk ratios while  $T^2$  for a synthesis of standardized mean differences will be in the metric of standardized mean differences. It would not be meaningful to compare the  $T^2$  values for two syntheses unless they were in the same metric. By contrast,  $I^2$  is on a ratio scale of 0 % to 100 %, and it is possible to compare this value from different syntheses.

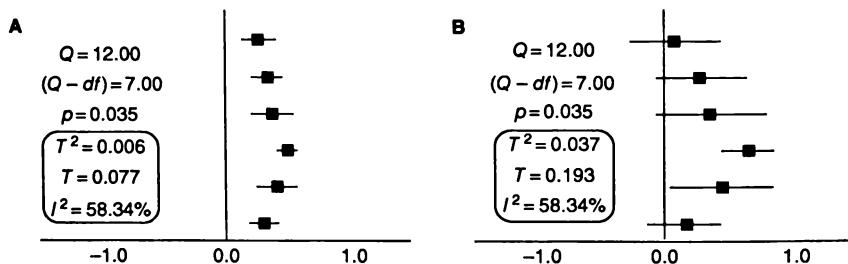


Figure 16.8 Factors affecting  $T^2$  but not  $I^2$ .

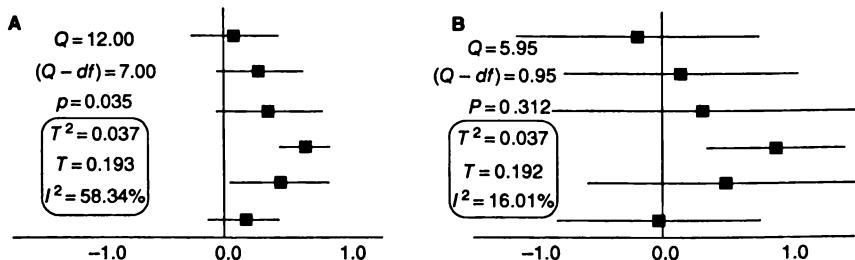


Figure 16.9 Factors affecting  $I^2$  but not  $T^2$ .

## In context

It is common for researchers who perform a meta-analysis to ask whether or not the effects are ‘heterogeneous’. As we have tried to show in this chapter, to report that the data are, or are not, heterogeneous, is not terribly informative. We need to consider what is meant by ‘heterogeneous’ and then respond with the relevant statistics.

Researchers often focus their attention on the  $Q$  statistic and its  $p$ -value. A significant  $Q$  conjures up images of effects that are widely dispersed and a nonsignificant  $Q$  is taken as assurance that the effects are consistent. This use of  $Q$  is clearly incorrect for two reasons. First, the  $Q$ -statistic and its  $p$ -value only address the viability of the null hypothesis (*Is the true dispersion exactly zero*) and not the amount of excess dispersion. Second,  $Q$  is sensitive to *relative* variance (the kind tracked by  $I^2$ ) and not absolute variance (the kind tracked by  $T^2$  and  $T$ ). These are two very different aspects of heterogeneity with very different implications (see below) and when a researcher reports  $Q$  alone, the consumer is not likely to make this distinction.

Researchers may also report  $T^2$ , since this is a key component in any random-effects meta-analysis. However, because  $T^2$  is a variance (and reported in the squared metric), it is not an intuitive measure. The problem of interpretation is even worse when the index is on a log scale.

One way to think about the substantive implications of the dispersion is to think about the range of effects, and how the utility of the intervention (or the importance of the relationship) varies over this range. The index that addresses this issue is  $T$ . Just as we might report the mean and standard deviation in a primary study to describe the distribution of scores, we can report the summary effect and standard deviation (that is,  $T$ ) in a meta-analysis to describe the distribution of true effects. For example, suppose that an effect size in the range of 0.0 to 0.20 would be considered trivial, 0.20 to 0.50 would be considered moderate, and 0.50 or above would be considered high. If the summary effect is 0.50 with a  $T$  of 0.10, most effects will fall in the range of 0.30 to 0.70, all in the moderate to high range. By contrast, if the summary effect is 0.50 with  $T$  of 0.20, some proportion of effects will fall in the trivial range. This approach works equally well when the index is on a log scale. Here, we compute the range of effects using this log scale, and then convert to the natural units for reporting. (Note that these intervals ignore uncertainty in the mean and in  $T$ . In Chapter 17 we show how to take account of this uncertainty when computing the prediction interval).

Another way to think about the substantive implications of the dispersion is to ask what proportion of the observed variance is real. Faced with a forest plot of risk ratios (for example) that range from 0.50 to 4.0, a researcher may be tempted to look for covariates that can explain the mechanism responsible for this dispersion. Before embarking on this quest, it makes sense to ask how much the true effects actually vary. If all true effects fall in a very narrow range, there may be no point in looking for moderators associated with this variance. By contrast, when there is substantial variation in true effects, it would be useful to look for these moderators. The prediction interval (Chapter 17) allows us to quantify the extent of heterogeneity on a meaningful scale.

Subgroup analyses and meta-regression may be employed to identify moderators that may explain some of this heterogeneity (Chapters 21 to 23).

An informative presentation of heterogeneity indices requires both a measure of the magnitude and a measure of uncertainty. Magnitude may be represented by the degree of true variation on the scale of the effect measure ( $T^2$ ) or the degree of inconsistency ( $I^2$ ), or both. Uncertainty over whether apparent heterogeneity is genuine may be expressed using the  $p$ -value for  $Q$  or using confidence intervals for  $T^2$  or  $I^2$ .

Note that uncertainty around  $T^2$  or  $I^2$  is often very large. If the studies themselves have poor precision (wide confidence intervals), this could mask the presence of real (possibly substantively important) heterogeneity, resulting in an estimate of zero for  $T^2$  and  $I^2$ . Therefore, it would be a mistake to interpret a  $T^2$  or  $I^2$  of zero as meaning that the effect sizes are consistent unless this is justified by confidence intervals for  $T^2$  and  $I^2$  that exclude large values.

## CONFIDENCE INTERVALS FOR $\tau^2$

If we assume that the effect sizes are normally distributed, the standard error of  $T^2$  may be estimated as follows.

First, compute

$$A = \left[ df + 2 \left( sw1 - \frac{sw2}{sw1} \right) \tau^2 + \left( sw2 - 2 \left( \frac{sw3}{sw1} \right) + \frac{(sw2)^2}{(sw1)^2} \right) \tau^4 \right], \quad (16.11)$$

where

$$\begin{aligned} sw1 &= \sum_{i=1}^k w_i, \\ sw2 &= \sum_{i=1}^n w_i^2, \end{aligned}$$

and

$$sw3 = \sum_{i=1}^n w_i^3.$$

Then, the variance of  $T^2$  is

$$V_{T^2} = 2 \times \left( \frac{A}{C^2} \right) \quad (16.12)$$

and its standard error is given by

$$SE_{T^2} = \sqrt{V_{T^2}}. \quad (16.13)$$

Because the distribution of  $T^2$  is not well approximated by a normal distribution, computing the confidence interval as the estimate of  $\tau^2$  plus or minus two standard errors will not yield very accurate confidence intervals unless the number of studies is very large. There are several methods for obtaining a confidence interval for  $\tau^2$ . A simple method is as follows.

First, if  $Q > (df + 1)$ , compute

$$B = 0.5 \times \frac{\ln(Q) - \ln(df)}{\sqrt{2Q} - \sqrt{2 \times df - 1}}, \quad (16.14)$$

or if  $Q \leq (df + 1)$ , compute

$$B = \sqrt{\frac{1}{2 \times (df - 1) \times \left(1 - \left(\frac{1}{3 \times (df - 1)^2}\right)\right)}}. \quad (16.15)$$

Then compute intermediate values

$$L = \text{Exp}\left(0.5 \times \ln\left(\frac{Q}{df}\right) - 1.96 \times B\right) \quad (16.16)$$

and

$$U = \text{Exp}\left(0.5 \times \ln\left(\frac{Q}{df}\right) + 1.96 \times B\right). \quad (16.17)$$

Finally, the 95% confidence intervals for  $\tau^2$  may then be obtained as

$$LL_{T^2} = \frac{df \times (L^2 - 1)}{C} \quad (16.18)$$

and

$$UL_{T^2} = \frac{df \times (U^2 - 1)}{C}. \quad (16.19)$$

Any value ( $T^2$ , lower limit or upper limit) that is computed as less than zero is set to zero. If the lower limit exceeds zero, then  $T^2$  should be statistically significant. However, since  $T^2$  is based on  $Q$ , and the sampling distribution of  $Q$  is better known, the preferred method would be to test  $Q$  for significance, and use this as the test for  $\tau^2$  being nonzero.

The 95% confidence interval for  $\tau$  may be obtained by taking the square roots of the confidence limits for  $\tau^2$ , namely

$$LL_T = \sqrt{LL_{T^2}}$$

and

$$UL_T = \sqrt{UL_{T^2}}.$$

## CONFIDENCE INTERVALS (OR UNCERTAINTY INTERVALS) FOR $I^2$

There are several methods for obtaining an interval to convey uncertainty in  $I^2$ . Because  $I^2$  does not estimate any underlying quantity, these intervals would be better described as uncertainty intervals rather than confidence intervals. However, we will continue to describe them as confidence intervals since the distinction is not practically important. A simple method to obtain confidence intervals is as follows.

First, if  $Q > (df + 1)$ , compute

$$B = 0.5 \times \frac{\ln(Q) - \ln(df)}{\sqrt{2Q} - \sqrt{2 \times df - 1}}, \quad (16.20)$$

or if  $Q \leq (df + 1)$ , compute

$$B = \sqrt{\frac{1}{2 \times (df - 1)} \left( 1 - \frac{1}{3 \times (df - 1)^2} \right)}. \quad (16.21)$$

Then

$$L = \exp \left( 0.5 \times \ln \left( \frac{Q}{df} \right) - 1.96 \times B \right) \quad (16.22)$$

and

$$U = \exp \left( 0.5 \times \ln \left( \frac{Q}{df} \right) + 1.96 \times B \right). \quad (16.23)$$

The 95% confidence intervals may then be obtained as

$$LL_{I^2} = \left( \frac{L^2 - 1}{L^2} \right) \times 100\% \quad (16.24)$$

and

$$UL_{I^2} = \left( \frac{U^2 - 1}{U^2} \right) \times 100\%. \quad (16.25)$$

Any value ( $I^2$ , lower limit or upper limit) that is computed as less than zero is set to zero.

If the lower limit of  $I^2$  exceeds zero, then  $I^2$  should be statistically significant. However, since  $I^2$  is based on  $Q$ , and the sampling distribution of  $Q$  is better known than the sampling distribution of  $I^2$ , the preferred method would be to test  $Q$  for significance, and use this as the test for  $I^2$  being nonzero.

Worked examples for all of these computations are included in Chapter 18.

### SUMMARY POINTS

- When we speak about dispersion in effect sizes from study to study we are usually concerned with the dispersion in true effect sizes, but the *observed* dispersion includes both true variance and random error.
- The mechanism used to isolate the true variance is to compare the observed dispersion with the amount we would expect to see if all studies shared a common effect size. The excess portion is assumed to reflect real differences among studies. This portion of the variance is then used to create several measures of heterogeneity.
- $Q$  is the weighted sum of squares (WSS) on a standardized scale. As a standard score it can be compared with the expected WSS (on the assumption that all studies share a common effect) to yield a test of the null hypothesis and also an estimate of the excess variance.
- $T^2$  is the variance of the true effects, on the same scale (squared) as the effects themselves. This value is used to assign study weights under the random-effects model.

- $T$  is the standard deviation of the true effects, on the same scale as the effects themselves. We can use this to estimate the distribution of true effects, and consider the substantive implications of this distribution.
- $I^2$  is the proportion of observed dispersion that is real, rather than spurious. It is not dependent on the scale, and is expressed as a ratio with a range of 0% to 100%.



# Prediction Intervals

---

### Introduction

Prediction intervals in primary studies

Prediction intervals in meta-analysis

Confidence intervals and prediction intervals

Comparing the confidence interval with the prediction interval

---

## INTRODUCTION

When we report the results of a meta-analysis we often focus on the summary effect size and its confidence interval. These give us an estimate of the mean effect size and its precision, *but they say nothing about how the true effects* are distributed about the summary effect.

In a fixed-effect analysis this is appropriate, since we assume that the true effect is the same in all studies. In a random-effects analysis, however, we need to consider not only the mean effect size, but also how the true effects are distributed about this mean. A mean effect size (say, a standardized mean difference) of 0.50 where all true effects are clustered in the range of 0.40 to 0.60 may have very different implications than the same mean effect where the true effects are scattered over the range of 0.00 to 1.00.

Our goal in this chapter is to show how we can use a prediction interval to describe the distribution of true effect sizes. We will review how the prediction interval is used in primary studies, and then show how the same mechanism can be used for meta-analysis.

## PREDICTION INTERVALS IN PRIMARY STUDIES

Suppose we are interested in math scores for a population of children. We want to create a prediction interval, defined as the interval within which a new student's score would fall if that student were selected at random from this population. The 80% prediction interval would include that score 80% of the time, the 95% interval would

include that score 95% of the time, and so on. As such, the interval yields an intuitive picture of the distribution of scores.

If we somehow knew the population mean ( $\mu$ ) and standard deviation ( $\sigma$ ), and were willing to assume that the scores are normally distributed, we could create a prediction interval, using

$$LL_{pred} = \mu - Z^\alpha \sqrt{\sigma^2} \quad (17.1)$$

and

$$UL_{pred} = \mu + Z^\alpha \sqrt{\sigma^2}, \quad (17.2)$$

where  $Z^\alpha$  is the  $Z$ -value corresponding to the desired confidence level (for the 95% interval,  $Z^\alpha$  would be 1.96). For example, if  $\mu$  is 0.50 and  $\sigma$  is 0.10, then the lower and upper limits of a 95% prediction interval are

$$LL_{pred} = 0.500 - 1.96 \times 0.100 = 0.3040$$

and

$$UL_{pred} = 0.500 + 1.96 \times 0.100 = 0.6960.$$

Formulas (17.1) and (17.2) are intuitive but are not useful in practice because they assume that we know both  $\mu$  and  $\sigma$  exactly. When these values are estimated from the sample (as they almost always are) we instead use the formulas

$$LL_{pred} = \bar{X} - t_{df}^\alpha \sqrt{S^2 + \frac{S^2}{n}} \quad (17.3)$$

and

$$UL_{pred} = \bar{X} + t_{df}^\alpha \sqrt{S^2 + \frac{S^2}{n}}, \quad (17.4)$$

where  $\bar{X}$  is the sample mean,  $t_{df}^\alpha$  is the  $t$ -value corresponding to (for example if  $\alpha = 0.05$ ) the 95% interval when there are  $df$  degrees of freedom, and  $S$  is the standard deviation of scores in the sample. These formulas have the same structure as (17.1) and (17.2) but to allow for error in the estimates of  $\mu$  and  $\sigma$  they incorporate the following changes. First, we multiply by  $t$  rather than  $Z$ . Second,  $t$  is multiplied by a quantity that involves both the variance of the observations (the standard deviation squared, or  $S^2$ ) and also the variance of the mean (the standard error squared, or  $S^2/n$ ).

For example, if  $X = 0.50$ ,  $S = 0.10$  and  $n = 30$ , then

$$LL_{pred} = 0.5000 - 2.0452 \times \sqrt{0.1000^2 + \frac{0.1000^2}{30}} = 0.2921$$

and

$$UL_{pred} = 0.5000 + 2.0452 \times \sqrt{0.1000^2 + \frac{0.1000^2}{30}} = 0.7079,$$

where a  $t$ -value of 2.045 corresponds to the  $t$ -value for alpha of 0.05 with 29  $df$ . In Excel, the function TINV(0.05,29) returns 2.0452.

Note that the prediction intervals based on the statistics (0.292 to 0.708) are wider than those based on the parameters (0.304 to 0.696).

## PREDICTION INTERVALS IN META-ANALYSIS

We can follow a similar approach in meta-analysis. If we somehow knew the mean effect size ( $\mu$ ) and the standard deviation of true effect sizes ( $\tau$ ), and were willing to assume that the effect sizes are normally distributed, we could create a prediction interval using

$$LL_{pred} = \mu - Z^\alpha \sqrt{\tau^2} \quad (17.5)$$

and

$$UL_{pred} = \mu + Z^\alpha \sqrt{\tau^2}. \quad (17.6)$$

This is similar to (17.1) and (17.2) except that  $\tau^2$ , the variance of the true effects in a meta-analysis, has replaced  $\sigma^2$ , the variance of the scores in a primary study.

Actually, we introduced this idea in Chapter 16 when we discussed the interpretation of  $T$  as an estimate of the standard deviation of true effect sizes. For example, if  $\mu$  is 0.358 and  $\tau^2$  is 0.0373, then the 95% prediction interval is

$$LL_{pred} = 0.358 - 1.96 \times \sqrt{0.0373} = -0.0205$$

and

$$UL_{pred} = 0.358 + 1.96 \times \sqrt{0.0373} = 0.7365.$$

In a forest plot we would typically use a simple line (from  $-0.020$  to  $0.737$ ) to represent the prediction interval, but in Figure 17.1 we use a bell curve to convey the idea that the true effect sizes are expected to be normally distributed within this range. Note that the bell curve has been truncated at either end ( $-0.020$  and  $0.737$ ) so that it covers 95% of expected true effects.

Formulas (17.5) and (17.6) assume that we actually know the values of  $\mu$  and  $\tau$ , and make no allowance for error in these estimates. Higgins *et al.* propose the following

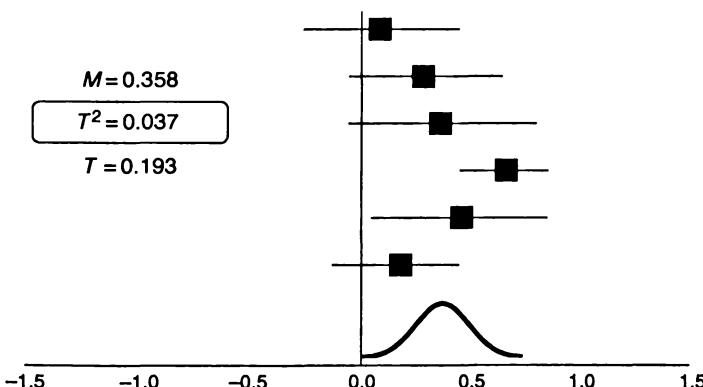


Figure 17.1 Prediction interval based on population parameters  $\mu$  and  $\tau^2$ .

formulas for computing a prediction interval when these values are estimated from the sample. The formulas are

$$LL_{pred} = M^* - t_{df}^\alpha \sqrt{T^2 + V_{M^*}} \quad (17.7)$$

and

$$UL_{pred} = M^* + t_{df}^\alpha \sqrt{T^2 + V_{M^*}}, \quad (17.8)$$

where  $M^*$  is the mean effect size in the sample,  $T^2$  is the sample estimate of the variance of true effect sizes, and  $V_{M^*}$  is the variance of  $M^*$ . The factor  $t$  is the  $t$ -value corresponding to (for example if  $\alpha = 0.05$ ) the 95% interval when there are  $df$  degrees of freedom.

These formulas have the same structure as (17.5) and (17.6) but we multiply by the  $t$ -value (rather than the Z-value) and apply this factor to a quantity that involves both the variance of the *true effects* ( $T^2$ ) and the variance of the *mean effect* ( $V_{M^*}$ ). The degrees of freedom ( $df$ ) is often taken as the number of studies minus 2 (that is,  $k - 2$ ).

For example, if  $k = 6$ ,  $M^* = 0.3582$ ,  $T^2 = 0.0373$ , and  $V_{M^*} = 0.0111$ , then

$$LL_{pred} = 0.3582 - 2.7764 \times \sqrt{0.0373 + 0.0111} = -0.2525$$

and

$$UL_{pred} = 0.3582 + 2.7764 \times \sqrt{0.0373 + 0.0111} = 0.9690.$$

The value 2.7764 is the  $t$ -value corresponding to alpha of 0.05 with 4  $df$ . In Excel, the function =TINV(0.05,4) returns 2.7764.

Figure 17.2 is identical to Figure 17.1 except that this time the prediction interval is based on the sample values  $M^*$  and  $T^2$  rather than the population parameters  $\mu$ , and  $\tau^2$ . Note that the bell curve is wider in Figure 17.2 (the 95% interval is  $-0.25$  to  $+0.97$ ) than in Figure 17.1 (where the interval was  $-0.02$  to  $+0.74$ ) which reflects the uncertainty in the estimates.

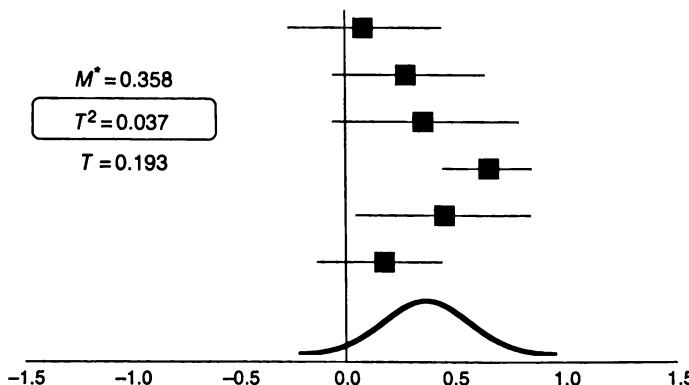


Figure 17.2 Prediction interval based on sample estimates  $M^*$  and  $T^2$ .

## CONFIDENCE INTERVALS AND PREDICTION INTERVALS

Traditionally, the summary line in a forest plot uses a diamond to depict the mean effect size (the center of the diamond) and its confidence interval (the width of the diamond). Now, we want to add a visual indicator of the prediction interval, and we do so by adding a horizontal line to either end of the diamond, as in Figure 17.3.

The meta-analysis line in the plot now shows *two distinct items of information*. First, in 95% of cases the mean effect size falls inside the diamond. Second, in 95% of cases the true effect in a new study will fall inside the horizontal lines. It is important to understand that these two items address two distinct issues. The confidence interval quantifies *the accuracy of the mean*, while the prediction interval addresses the actual *dispersion of effect sizes*, and the two measures are not interchangeable.

As always, how we choose to interpret the effects depends on our goals. We may want to focus on the null effect. If the full diamond exceeds zero then we are 95% certain that the mean effect size exceeds zero. If the full prediction interval exceeds zero then the true effect in 95% of new studies will exceed zero.

Or, we may want to focus on a clinically important effect (say, a standardized mean difference of 0.20). If the full diamond exceeds 0.20 then we are 95% certain that the mean effect size exceeds 0.20. If the full prediction interval exceeds 0.20 then the true effect in 95% of new studies will exceed 0.20.

## COMPARING THE CONFIDENCE INTERVAL WITH THE PREDICTION INTERVAL

Earlier, we showed the computation of the prediction interval for a meta-analysis with six studies. Suppose that the meta-analysis included more studies (24, 60, or 1002) with the same pattern as in the first six. In other words, we have the same within-study

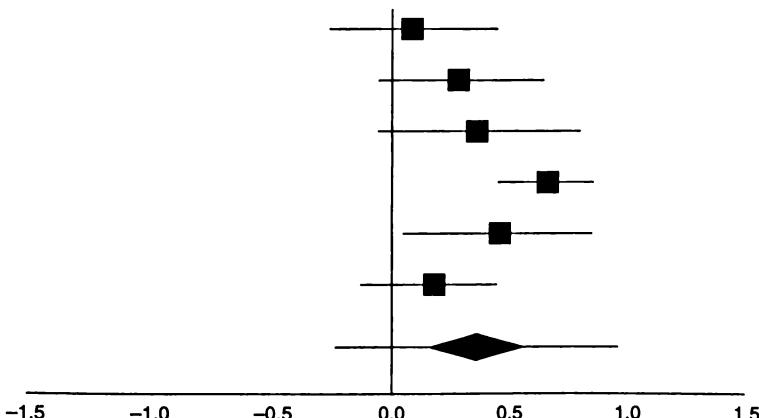
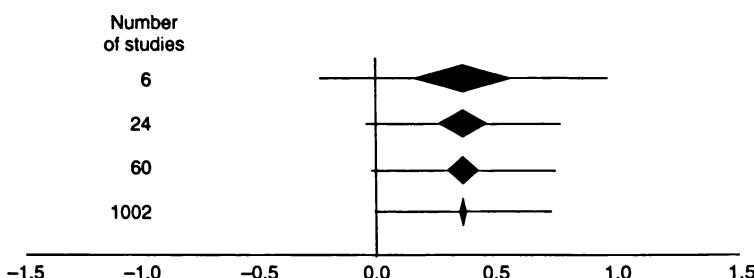


Figure 17.3 Simultaneous display of confidence interval and prediction interval.



**Figure 17.4** Impact of number of studies on confidence interval and prediction interval.

error and the same pattern of dispersion but a more precise estimate of the mean effect size and of the true between-studies dispersion.

In Figure 17.4 we illustrate what the confidence interval and the prediction interval would be for these four hypothetical analyses. While the specific pattern shown here is unique to this analysis, the general trend will apply to any analysis.

With six studies the confidence interval (the diamond) is quite wide, but with 60 studies its width is cut by about half, and with 1002 studies its width is trivial. This follows from the formula for a confidence interval, which is

$$CI_{M^*} = M^* \pm Z \sqrt{V_{M^*}} \quad (17.9)$$

The confidence interval reflects only error ( $V_{M^*}$ ), and so we see a consistent decline in the confidence interval width as the number of studies goes from 6 to 1002. With an infinite number of studies the error would approach zero, and so the width of the confidence interval would approach zero.

By contrast, the width of the prediction interval (the line) drops sharply as the number of studies goes from 6 to 24 but shows almost no change beyond that point. This follows from the formula for a prediction interval, which is

$$PI = M^* \pm t \sqrt{T^2 + V_{M^*}} \quad (17.10)$$

The interval is based on error in estimating the mean ( $V_{M^*}$ ), which is dependent on the number of studies. The interval is based also on the variance of the studies,  $T^2$ , which is not affected by the number of studies. In this example, as the number of studies increases from 6 to 24,  $V_{M^*}$  decreases and therefore the interval narrows. Beyond that point the decrease in  $V_{M^*}$  is trivial (and  $T^2$  remains constant), so the prediction interval shows little change. With an infinite number of studies, the interval would approach  $\mu$  plus/minus 1.96  $\tau$ .

## SUMMARY POINTS

- For a random-effects analysis we want to know both the mean effect size and also how the true effects are distributed about the mean.
- The precision of the mean is addressed by the confidence interval. Since the confidence interval reflects only error of estimation of the mean, with an infinite number of studies its width would approach zero.
- The distribution of true effect sizes is addressed by the prediction interval. Since the prediction interval incorporates true dispersion as well as error, with an infinite number of studies it will approach the actual dispersion of true effect sizes.
- The summary effect in a forest plot has traditionally been represented by a diamond which corresponds to the confidence interval. For random-effects analyses we can modify this to display both the confidence interval and the prediction interval.

## Further Reading

- Borenstein, M., Higgins J.P.T., Hedges, L.V., and Rothstein, H.R. (2017). Basics of meta-analysis:  $I^2$  is not an absolute measure of heterogeneity. *Research Synthesis Methods* 8 (1): 5–18. <https://doi.org/10.1002/jrsm.1230>.
- Borenstein, M. (2019). *Common Mistakes in Meta-Analysis and How to Avoid Them*. Englewood, NJ: Biostat, Inc.
- Borenstein, M. (2020). Research Note: In a meta-analysis, the  $I^2$  index does not tell us how much the effect size varies across studies. *Journal of Physiotherapy* 66 (2): 135–139. <https://doi.org/10.1016/j.jphys.2020.02.011>.



# Worked Examples (Part 2)

---

### Introduction

Worked example for continuous data (Part 2)

Worked example for binary data (Part 2)

Worked example for correlational data (Part 2)

---

### INTRODUCTION

In Chapter 14 we presented worked examples for computing a summary effect using continuous, binary, and correlational data. Here, we continue with the same three data sets and show how to compute the measures of heterogeneity discussed in Chapters 16 and 17.

These computations are also included in Excel spreadsheets that can be downloaded from the book's website

### WORKED EXAMPLE FOR CONTINUOUS DATA (PART 2)

On page 81 we showed how to compute the effect size and variance for each study. Here, we proceed from that point.

Using results in Table 18.1, the summary effect is given by

$$M = \frac{101.171}{244.215} = 0.4143,$$

which value is used in the column labeled *Mean* in Table 18.2.

Then, using (16.1) we sum the values in the final column of Table 18.2,

$$Q = \sum_{i=1}^k W_i(Y_i - M)^2 = 12.0033.$$

Or, using (12.3) and results in Table 18.1,

$$Q = 53.915 - \frac{(101.171)^2}{244.215} = 12.0033.$$

**Table 18.1** Dataset 1 – Part D (intermediate computations).

Study	Effect $Y$	Variance $V_Y$	Weight $W$	Calculated quantities			
				$WY$	$WY^2$	$W^2$	$W^3$
Carroll	0.095	0.033	30.352	2.869	0.271	921.21	27960.25
Grant	0.277	0.031	32.568	9.033	2.505	1060.68	34544.41
Peck	0.367	0.050	20.048	7.349	2.694	401.93	8058.00
Donat	0.664	0.011	95.111	63.190	41.983	9046.01	860371.10
Stewart	0.462	0.043	23.439	10.824	4.999	549.37	12876.47
Young	0.185	0.023	42.698	7.906	1.464	1823.12	77843.29
Sum			244.215	101.171	53.915	13802.33	1021653.52

**Table 18.2** Dataset 1 – Part E (variance computations).

Study	Effect $Y$	Variance $V_Y$	Weight $W$	Mean $M$	Calculated quantities	
					$(Y - M)^2$	$W(Y - M)^2$
Carroll	0.095	0.033	30.352	0.414	0.102	3.103
Grant	0.277	0.031	32.568	0.414	0.019	0.610
Peck	0.367	0.050	20.048	0.414	0.002	0.046
Donat	0.664	0.011	95.111	0.414	0.063	5.950
Stewart	0.462	0.043	23.439	0.414	0.002	0.053
Young	0.185	0.023	42.698	0.414	0.052	2.241
Sum						12.003

Under the assumption that all studies share a common effect, the expected value of  $Q$  is given by

$$df = 6 - 1 = 5$$

where  $k$  is the number of studies. The difference,

$$12.003 - 5 = 7.003,$$

is the excess value which we attribute to differences in the true effect sizes.

The  $p$ -value for  $Q = 12.003$  with  $df = 5$ , is 0.035. In Excel, the function =CHIDIST(12.003,5) returns 0.035. If we are using 0.10 or 0.05 as the criterion for statistical significance, we would reject the null hypothesis that all the studies share a common effect size, and accept the alternative, that the true effect is not the same in all studies.

Then, using formulas (16.6), (16.5), (16.8), and (16.9),

$$C = 244.215 - \left( \frac{13802.33}{244.215} \right) = 187.6978,$$

$$T^2 = \frac{12.003 - 5}{187.698} = 0.0373,$$

$$T = \sqrt{0.0373} = 0.1932,$$

and

$$I^2 = \left( \frac{12.003 - 5}{12.003} \right) \times 100\% = 58.34\%.$$

To compute the standard error of  $T^2$  (from (16.11) to (16.13)), we have  $sw1 = 244.215$ ,  $sw2 = 13,802.33$ , and  $sw3 = 1,021,653.52$ , so that

$$\begin{aligned} A &= \left[ df + 2 \left( 244 - \frac{13802}{244} \right) 0.0373 \right. \\ &\quad \left. + \left( 13802 - 2 \left( \frac{1021653}{244} \right) + \frac{(13802)^2}{(244)^2} \right) 0.0373^2 \right] = 31.0202. \end{aligned}$$

Then, the variance of  $T^2$  is

$$V_{T^2} = 2 \times \left( \frac{31.020}{187.698^2} \right) = 0.0018,$$

and its standard error is given by

$$SE_{T^2} = \sqrt{0.0018} = 0.0420.$$

Since  $Q = 12.003 > 6 = (df + 1)$ , we compute, from (16.14) to (16.19),

$$B = 0.5 \times \frac{\ln(12.0033) - \ln(5)}{\sqrt{2 \times 12.0033} - \sqrt{2 \times 5 - 1}} = 0.2305.$$

Then compute intermediate values

$$L = \text{Exp} \left( 0.5 \times \ln \left( \frac{12.003}{5} \right) - 1.96 \times 0.2305 \right) = 0.9862$$

and

$$U = \text{Exp} \left( 0.5 \times \ln \left( \frac{12.003}{5} \right) + 1.96 \times 0.2305 \right) = 2.4343.$$

Finally, the 95% confidence intervals for  $\tau^2$  may be obtained as

$$LL_{T^2} = \frac{5 \times (0.9862^2 - 1)}{187.698} = -0.0007,$$

which is set to zero, and

$$UL_{T^2} = \frac{5 \times (2.4343^2 - 1)}{187.698} = 0.1312.$$

The 95% confidence interval for  $\tau$  may be obtained by taking the square roots of the confidence limits for  $\tau^2$ , namely

$$LL_T = \sqrt{0.0} = 0.0,$$

and

$$UL_T = \sqrt{0.1312} = 0.3622.$$

### Confidence intervals for $I^2$

Since  $12.003 > (5 + 1)$  we compute, using formulas (16.20) through (16.25),

$$B = 0.5 \times \frac{\ln(12.003) - \ln(5)}{\sqrt{2 \times 12.003} - \sqrt{2 \times 5 - 1}} = 0.2305.$$

Compute intermediate values

$$L = \exp \left( 0.5 \times \ln \left( \frac{12.003}{5} \right) - 1.96 \times 0.2305 \right) = 0.9862$$

and

$$U = \exp \left( 0.5 \times \ln \left( \frac{12.003}{5} \right) + 1.96 \times 0.2305 \right) = 2.4343.$$

The 95% confidence intervals may then be obtained as

$$LL_{P2} = \left( \frac{0.9862^2 - 1}{0.9862^2} \right) \times 100\% = -2.82\%,$$

which is set to zero, and

$$UL_{P2} = \left( \frac{2.4343^2 - 1}{2.4343^2} \right) \times 100\% = 83.12\%.$$

To obtain a 95% prediction interval for the true standardized mean difference in a future study, we use the random-effects weighted mean and its variance computed in (14.1) and (14.2),  $M^* = 0.3582$  and  $V_{M^*} = 0.0111$  and compute, from (17.7) and (17.8),

$$t_{4}^{0.05} = 2.7764,$$

$$LL_{pred} = 0.3582 - 2.7764 \times \sqrt{0.0373 + 0.0111} = -0.2525,$$

and

$$UL_{pred} = 0.3582 + 2.7764 \times \sqrt{0.0373 + 0.0111} = 0.9690.$$

This prediction interval is plotted in Figure 18.1.

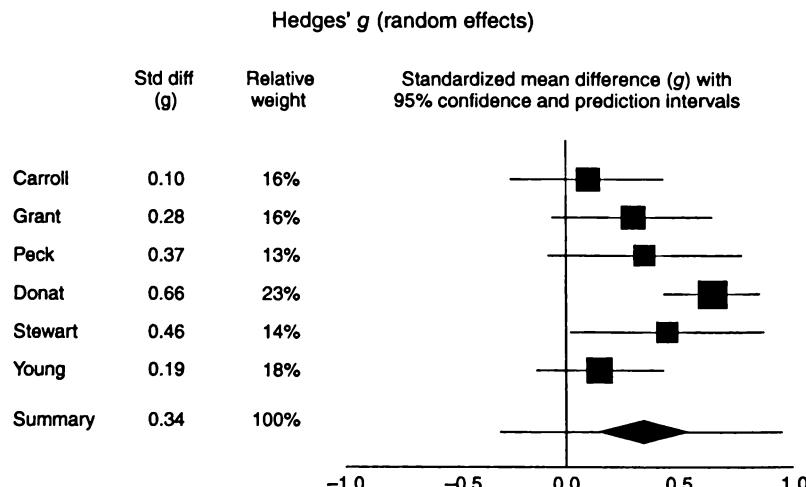


Figure 18.1 Forest plot of Dataset 1 – random-effects weights with prediction interval.

### WORKED EXAMPLE FOR BINARY DATA (PART 2)

On page 85 we showed how to compute the effect size (here, the log odds ratio) and variance for each study. Here, we proceed from that point.

Using results in Table 18.3, the summary effect is given by

$$M = \frac{-30.594}{42.248} = -0.7241,$$

which value is used in the column labeled *Mean* in Table 18.4.

Then, using (16.1) we sum the values in the final column of Table 18.4,

$$Q = \sum_{i=1}^k W_i(Y_i - M)^2 = 10.5512.$$

Or, using (12.3) and results in Table 18.3,

$$Q = 32.705 - \frac{(-30.594)^2}{42.248} = 10.5512.$$

Under the assumption that all studies share a common effect, the expected value of  $Q$  is given by

$$df = 6 - 1 = 5,$$

**Table 18.3** Dataset 2 – Part D (intermediate computations).

Study	Effect $Y$	Variance $V_Y$	Weight $W$	Calculated quantities			
				$WY$	$WY^2$	$W^2$	$W^3$
Saint	-0.366	0.185	5.402	-1.978	0.724	29.18	157.66
Kelly	-0.288	0.290	3.453	-0.993	0.286	11.92	41.18
Pilbeam	-0.384	0.156	6.427	-2.469	0.948	41.30	265.42
Lane	-1.322	0.058	17.155	-22.675	29.971	294.30	5048.71
Wright	-0.417	0.282	3.551	-1.480	0.617	12.61	44.76
Day	-0.159	0.160	6.260	-0.998	0.159	39.19	245.33
Sum			42.248	-30.594	32.705	428.50	5803.06

**Table 18.4** Dataset 2 – Part E (variance computations).

Study	Effect $Y$	Variance $V_Y$	Weight $W$	Mean $M$	Calculated quantities	
					$(Y - M)^2$	$W(Y - M)^2$
Saint	-0.366	0.185	5.402	-0.724	0.128	0.692
Kelly	-0.288	0.290	3.453	-0.724	0.191	0.658
Pilbeam	-0.384	0.156	6.427	-0.724	0.116	0.743
Lane	-1.322	0.058	17.155	-0.724	0.357	6.127
Wright	-0.417	0.282	3.551	-0.724	0.094	0.335
Day	-0.159	0.160	6.260	-0.724	0.319	1.996
Sum						10.551

where  $k$  is the number of studies. The difference,

$$10.5512 - 5 = 5.5512$$

is the excess value which we attribute to differences in the true effect sizes.

The  $p$ -value for  $Q = 10.551$  with  $df = 5$ , is 0.0610. In Excel, the function =CHIDIST(10.551,5) returns 0.0610. If we are using 0.10 as the criterion for statistical significance, we would reject the null hypothesis that all the studies share a common effect size, and accept the alternative, that the true effect is not the same in all studies. If we are using 0.05 as the criterion, we would not have sufficient evidence to reject the null hypothesis (but would not conclude that the effects are homogeneous, since the nonsignificant  $p$ -value could be due to inadequate statistical power).

Then, using formulas (16.6), (16.5), (16.8), and (16.9),

$$C = 42.248 - \left( \frac{428.50}{42.248} \right) = 32.1052,$$

$$T^2 = \frac{10.5512 - 5}{32.1052} = 0.1729,$$

$$T = \sqrt{0.1729} = 0.4158,$$

and

$$I^2 = \left( \frac{10.5512 - 5}{10.5512} \right) \times 100 = 52.61\%.$$

To compute the standard error of  $T^2$  (from (16.11) to (16.13)), we have  $sw1 = 42.25$ ,  $sw2 = 428.5$ , and  $sw3 = 5,803.1$ , so that

$$A = \left[ df + 2 \left( 42.25 - \frac{428.5}{42.25} \right) 0.1729 + \left( 428.5 - 2 \left( \frac{5803.1}{42.25} \right) + \frac{(428.5)^2}{(42.25)^2} \right) 0.1729^2 \right] = 23.7754.$$

Then, the variance of  $T^2$  is

$$V_{T^2} = 2 \times \left( \frac{23.7754}{32.1052^2} \right) = 0.0461$$

and its standard error is given by

$$SE_{T^2} = \sqrt{0.0461} = 0.2148.$$

Since  $Q = 10.5512 > 6.5 (df + 1)$ , we compute, from (16.14) to (16.19),

$$B = 0.5 \times \frac{\ln(10.5512) - \ln(5)}{\sqrt{2 \times 10.5512} - \sqrt{2 \times 5 - 1}} = 0.2343.$$

Then compute intermediate values

$$L = \text{Exp} \left( 0.5 \times \ln \left( \frac{10.5512}{5} \right) - 1.96 \times 0.2343 \right) = 0.9178$$

and

$$U = \text{Exp} \left( 0.5 \times \ln \left( \frac{10.5512}{5} \right) + 1.96 \times 0.2343 \right) = 2.2993.$$

Finally, the 95% confidence intervals for  $\tau^2$  may then be obtained as

$$LL_{\tau^2} = \frac{5 \times (0.9178^2 - 1)}{32.1052} = -0.0246,$$

which is set to zero, and

$$UL_{\tau^2} = \frac{5 \times (2.2993^2 - 1)}{32.1052} = 0.6676.$$

The 95% confidence interval for  $\tau$  may be obtained by taking the square roots of the confidence limits for  $\tau^2$ , namely

$$LL_{\tau} = \sqrt{0.0} = 0.0,$$

and

$$UL_{\tau} = \sqrt{0.6676} = 0.8171.$$

### Confidence intervals for $I^2$

Since  $10.5512 > (5 + 1)$  we compute, using formulas (16.20) through (16.25),

$$B = 0.5 \times \frac{\ln(10.5512) - \ln(5)}{\sqrt{2 \times 10.5512} - \sqrt{2 \times 5} - 1} = 0.2343,$$

then compute intermediate values

$$L = \text{Exp} \left( 0.5 \times \ln \left( \frac{10.5512}{5} \right) - 1.96 \times 0.2343 \right) = 0.9178$$

and

$$U = \text{Exp} \left( 0.5 \times \ln \left( \frac{10.5512}{5} \right) + 1.96 \times 0.2343 \right) = 2.2993.$$

The 95% confidence intervals may then be obtained as

$$LL_{I^2} = \left( \frac{0.9178^2 - 1}{0.9178^2} \right) \times 100\% = -18.72\%,$$

which is set to zero, and

$$UL_{I^2} = \left( \frac{2.2993^2 - 1}{2.2993^2} \right) \times 100\% = 81.09\%.$$

To obtain a 95% prediction interval for the true log odds ratio in a future study, we use the random-effects weighted mean and its variance computed in (14.3) and (14.4),  $M^* = -0.5663$  and  $V_{M^*} = 0.0570$ , and compute, from (17.7) and (17.8),

$$t_4^{0.05} = 2.7764,$$

$$LL_{pred} = -0.5663 - 2.7764 \times \sqrt{0.1729 + 0.0570} = -1.8977,$$

and

$$UL_{pred} = -0.5663 + 2.7764 \times \sqrt{0.1729 + 0.0570} = 0.7651.$$

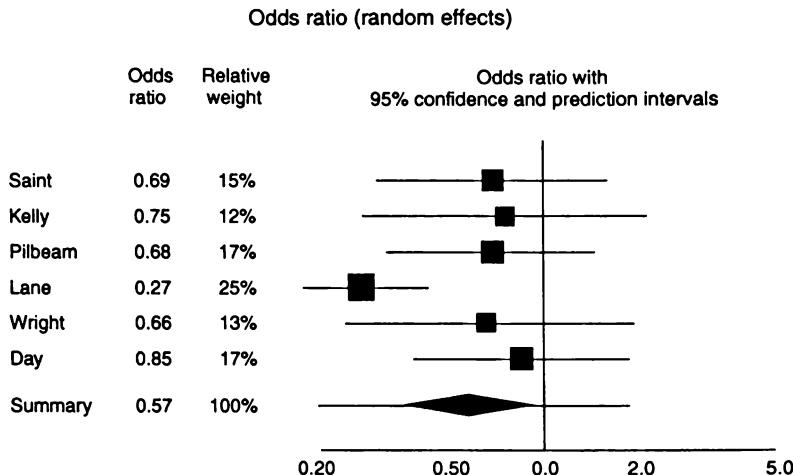


Figure 18.2 Forest plot of Dataset 2 – random-effects weights with prediction interval.

These limits are computed on a log scale. We can convert the limits to the odds ratio scale using

$$LL_{pred} = \exp(-1.8977) = 0.1499$$

and

$$UL_{pred} = \exp(0.7651) = 2.1492.$$

This prediction interval is plotted in Figure 18.2.

### WORKED EXAMPLE FOR CORRELATIONAL DATA (PART 2)

On page 90 we showed how to compute the effect size (here, the Fisher's  $z$  transformation of the correlation coefficient) and variance for each study. Here, we proceed from that point.

Using results in Table 18.5, the summary effect is given by

$$M = \frac{242.650}{647.000} = 0.3750,$$

which value is used in the column labeled *Mean* in Table 18.6.

Then, using (16.1) we sum the values in the final column of Table 18.6,

$$Q = \sum_{i=1}^k W_i(Y_i - M)^2 = 36.1437.$$

Or, using (12.3) and results in Table 18.5,

$$Q = 127.147 - \frac{(242.650)^2}{647.000} = 36.1437.$$

**Table 18.5** Dataset 3 – Part D (intermediate computations).

Study	Effect Y	Variance $V_Y$	Weight W	Calculated quantities			
				WY	WY <sup>2</sup>	W <sup>2</sup>	W <sup>3</sup>
Fonda	0.549	0.027	37.000	20.324	11.164	1369.00	50653.00
Newman	0.693	0.011	87.000	60.304	41.799	7569.00	658503.00
Grant	0.424	0.045	22.000	9.320	3.949	484.00	10648.00
Granger	0.203	0.003	397.000	80.485	16.317	157609.00	62570773.00
Milland	0.867	0.018	57.000	49.436	42.876	3249.00	185193.00
Finch	0.485	0.021	47.000	22.781	11.042	2209.00	103823.00
Sum			647.000	242.650	127.147	172489.00	63579593.00

**Table 18.6** Dataset 3 – Part E (variance computations).

Study	Effect Y	Variance $V_Y$	Weight W	Mean M	Calculated quantities	
					(Y - M) <sup>2</sup>	W(Y - M) <sup>2</sup>
Fonda	0.549	0.027	37.000	0.375	0.030	1.124
Newman	0.693	0.011	87.000	0.375	0.101	8.804
Grant	0.424	0.045	22.000	0.375	0.002	0.052
Granger	0.203	0.003	397.000	0.375	0.030	11.787
Milland	0.867	0.018	57.000	0.375	0.242	13.812
Finch	0.485	0.021	47.000	0.375	0.012	0.565
Sum						36.144

Under the assumption that all studies share a common effect, the expected value of  $Q$  is given by

$$df = 6 - 1 = 5,$$

where  $k$  is the number of studies. The difference,

$$36.1437 - 5 = 31.1437,$$

is the excess value which we attribute to differences in the true effect sizes.

The  $p$ -value for  $Q = 36.1437$  with  $df = 5$ , is less than 0.0001. In Excel, the function =CHIDIST(36.1437,5) returns < 0.0001. If we are using 0.10 or 0.05 as the criterion for statistical significance, we would reject the null hypothesis that all the studies share a common effect size, and accept the alternative, that the true effect is not the same in all studies.

Then, using formulas (16.6), (16.5), (16.8), and (16.9),

$$C = 647.000 - \left( \frac{172489.00}{647.000} \right) = 380.4019,$$

$$T^2 = \frac{36.1437 - 5}{380.4019} = 0.0819,$$

$$T = \sqrt{0.0819} = 0.28613,$$

and

$$I^2 = \left( \frac{36.1437 - 5}{36.1437} \right) \times 100 = 86.17\%.$$

To compute the standard error of  $T^2$  (from (16.11) to (16.13)), we have  $sw1 = 647.00$ ,  $sw2 = 172489.00$ , and  $sw3 = 63,579,593.00$ , so that

$$A = \left[ df + 2 \left( 647.00 - \frac{172489}{647.00} \right) 0.0819 + \left( 172489 - 2 \left( \frac{63579593}{647.00} \right) + \frac{(172489)^2}{(647.00)^2} \right) 0.0819^2 \right] = 382.4983.$$

Then, the variance of  $T^2$  is

$$V_{T^2} = 2 \times \left( \frac{382.4983}{380.4019^2} \right) = 0.0053,$$

and its standard error is given by

$$SE_{T^2} = \sqrt{0.0053} = 0.0727.$$

Since  $Q = 36.1437 > 6 = (df + 1)$ , we compute, from (16.14) to (16.19),

$$B = 0.5 \times \frac{\ln(36.1437) - \ln(5)}{\sqrt{2 \times 36.1437} - \sqrt{2 \times 5 - 1}} = 0.1798.$$

Then compute intermediate values

$$L = \text{Exp} \left( 0.5 \times \ln \left( \frac{36.1437}{5} \right) - 1.96 \times 0.1798 \right) = 1.8903$$

and

$$U = \text{Exp} \left( 0.5 \times \ln \left( \frac{36.1437}{5} \right) + 1.96 \times 0.1798 \right) = 3.8242.$$

Finally, the 95% confidence intervals for  $\tau^2$  may then be obtained as

$$LL_{T^2} = \frac{5 \times (1.890^2 - 1)}{380.4019} = 0.0338$$

and

$$UL_{T^2} = \frac{5 \times (3.8242^2 - 1)}{380.4019} = 0.1791.$$

The 95% confidence interval for  $\tau$  may be obtained by taking the square roots of the confidence limits for  $\tau^2$ , namely

$$LL_{\tau} = \sqrt{0.0338} = 0.1839,$$

and

$$UL_{\tau} = \sqrt{0.1791} = 0.4232.$$

### Confidence intervals for $I^2$

Since  $Q = 36.1437 > 6 = (df + 1)$ , we compute, from (16.20),

$$B = 0.5 \times \frac{\ln(36.1437) - \ln(5)}{\sqrt{2 \times 36.1437} - \sqrt{2 \times 5 - 1}} = 0.1798,$$

then compute intermediate values

$$L = \text{Exp} \left( 0.5 \times \ln \left( \frac{36.1437}{5} \right) - 1.96 \times 0.1798 \right) = 1.8903$$

and

$$U = \text{Exp} \left( 0.5 \times \ln \left( \frac{36.1437}{5} \right) + 1.96 \times 0.1798 \right) = 3.8242.$$

The 95% confidence intervals may then be obtained as

$$LL_{12} = \left( \frac{1.8903^2 - 1}{1.8903^2} \right) \times 100\% = 72.01\%,$$

and

$$UL_{12} = \left( \frac{3.8241^2 - 1}{3.8241^2} \right) \times 100\% = 93.16\%.$$

To obtain a 95% prediction interval for the true Fisher's  $z$  in a future study, we use the random-effects weighted mean and its variance computed in (14.5) and (14.6),  $M^* = 0.5328$  and  $V_{M^*} = 0.0168$  and compute, from (17.7) and (17.8),

$$t_4^{0.05} = 2.7764,$$

$$LL_{pred} = 0.5328 - 2.7764 \times \sqrt{0.0819 + 0.0168} = -0.3396,$$

and

$$UL_{pred} = 0.5328 + 2.7764 \times \sqrt{0.0819 + 0.0168} = 1.4051.$$

These limits are in the Fisher's  $z$  metric. We can convert the limits to the correlation scale using

$$LL_{pred} = \frac{e^{(2x-0.3396)} - 1}{e^{(2x-0.3396)} + 1} = -0.3271$$

and

$$UL_{pred} = \frac{e^{(2x1.4051)} - 1}{e^{(2x1.4051)} + 1} = 0.8865.$$

This prediction interval is plotted in Figure 18.3.

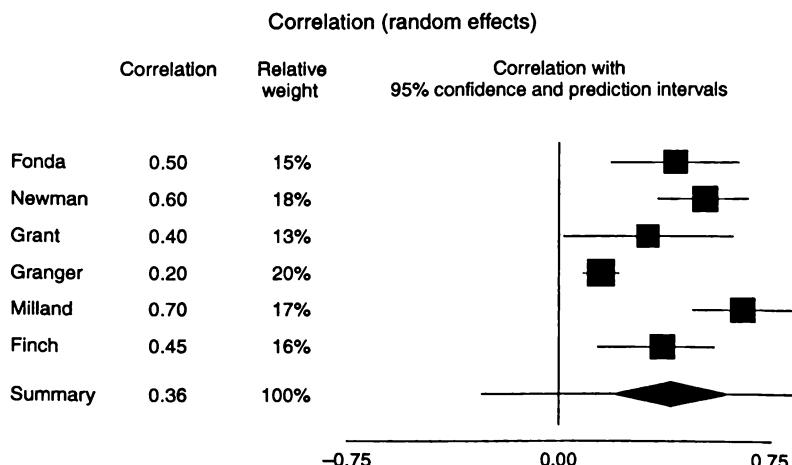


Figure 18.3 Forest plot of Dataset 3 – random-effects weights with prediction interval.

**SUMMARY POINTS**

- This chapter includes worked examples showing how to compute the summary effect using fixed-effect and random-effects models.
- For the standardized mean difference we work with the effect sizes directly.
- For ratios we work with the log transformed data.
- For correlations we work with the Fisher's  $z$  transformed data.
- These worked examples are available as Excel files on the book's website ([www.Introduction-to-Meta-Analysis.com](http://www.Introduction-to-Meta-Analysis.com)).

# An Intuitive Look at Heterogeneity

---

Introduction

Motivating example

The  $Q$ -value and the  $p$ -value do not tell us how much the effect size varies

The confidence interval does not tell us how much the effect size varies

The  $I^2$  statistic does not tell us how much the effect size varies

What  $I^2$  tells us

The  $I^2$  index vs. the prediction interval

The prediction interval

Prediction interval is clear, concise, and relevant

Computing the prediction interval

How to use  $I^2$

How to explain heterogeneity

How much does the effect size vary across studies?

Caveats

Conclusion

Further reading

The meaning of  $I^2$  in Figure 19.2

---

## INTRODUCTION

In previous chapters, we explained the meaning of the various indices employed to quantify heterogeneity. While many researchers understand the distinction between these indices in the abstract, relatively few actually put this knowledge into practice. Our goal in this chapter is to provide practical advice about how to think about heterogeneity.

The potential utility of an intervention depends not only on the mean effect size, but also on the dispersion of effects about that mean. We need to know if the intervention has essentially the same impact in all populations; or has a trivial impact in some populations and a large impact in others; or if it is harmful in some populations and helpful in others. When researchers ask about heterogeneity, they are asking which of these descriptions applies. However, the statistics typically reported

for heterogeneity ( $Q$ ,  $T^2$ ,  $I^2$ ) do not directly address this question. In this chapter, we highlight the prediction interval, the statistic that reports the range of true effects. This statistic provides the information that we need, and that many think is being provided by the other statistics.

While it is important to report the relevant statistics, it is also imperative to understand the limitations of these statistics. We need a reasonable number of studies to yield reliable estimates of any statistics related to heterogeneity. This applies to  $I^2$  and  $T^2$  as well as to the prediction interval. An estimate of heterogeneity that is based on a handful of studies (or fewer) is not likely to be reliable.

### MOTIVATING EXAMPLE

Ronksley, Brien, Turner, Mukamal, and Ghali (2011) published a meta-analysis in *BMJ* that looked at the relationship between alcohol consumption and all-cause mortality. The mean risk ratio was 0.87, which tells us that persons classified as drinkers had a lower risk of death than those classified as nondrinkers. The confidence interval is 0.83 to 0.91, and the Z-value for a test of the null hypothesis is 5.77 with a corresponding *p*-value of <0.001.

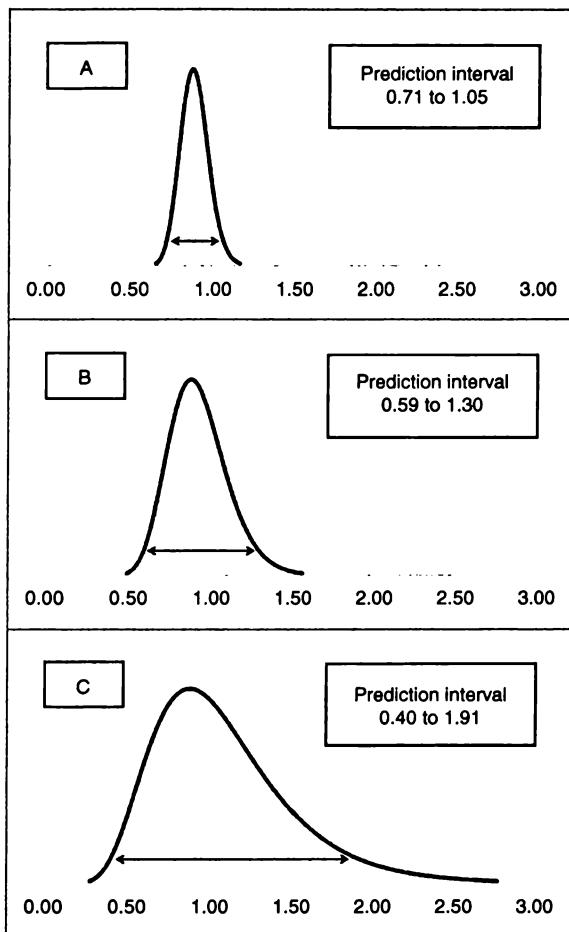
On this basis, we conclude that the mortality risk is 13% lower for drinkers, *on average*. However, we still need to address heterogeneity. That is, we need to know if the distribution of effects resembles panel A, B, or C in Figure 19.1. In each panel, the distribution may be summarized by means of the prediction interval, denoted by an arrow. The true effect size in 95% of all populations will fall inside that interval.

If the distribution of effects resembles panel A, we might report that the drinkers have a lower risk of death than nondrinkers in virtually all populations. At one extreme, there are some populations where the risk of death is 29% *lower* for drinkers. At the other extreme, there are a few populations where the risk of death is 5% *higher* for drinkers. As such, the relationship between drinking and mortality is relatively modest but also relatively consistent.

If the distribution of effects resembles panel B, we might report that the drinkers have a lower risk *on average*, but there is substantial variation in this relationship. At one extreme, there are some populations where the risk of death is 41% *lower* for drinkers. At the other extreme, there are some populations where the risk of death is 30% *higher* for drinkers.

Finally, if the distribution of effects resembles panel C, we might report that there is so much variation in the effect that the mean effect is of little relevance. At one extreme, there are some populations where the risk of death is 60% *lower* for drinkers. At the other extreme, there are some populations where the risk of death is 91% *higher* for drinkers.

These interpretations of the numbers are subjective, and others will characterize the implications of the heterogeneity differently. That discussion is necessary and welcome. However, to have an informed discussion about the implications of the dispersion, we must first know if the distribution resembles panel A, B, or C.



**Figure 19.1** Alcohol use and mortality. Risk ratio  $< 1$  favors drinkers. Three possible distributions of true effects.

Note that the distribution of effects is assumed to be symmetric in log units. It appears to be skewed because the plot uses the risk ratio rather than the log risk ratio on the X-axis.

#### THE Q-VALUE AND THE p-VALUE DO NOT TELL US HOW MUCH THE EFFECT SIZE VARIES

The statistics that most papers report for heterogeneity include the *Q*-value and the *p*-value. In the current analysis, the *Q*-value is 96.85 with 32 degrees of freedom, and

the  $p$ -value for a test of the null hypothesis (that the true effect size is the same in all studies) is  $<0.001$ . Based on these statistics, there is no way of knowing whether the distribution of effects resembles A, B, or C. The  $Q$ -value is the sum of squared deviations on a standardized scale and is driven by the number of studies and the extent of dispersion. The same applies to the  $p$ -value. Therefore, neither of these can serve as a surrogate for the amount of dispersion.

### THE CONFIDENCE INTERVAL DOES NOT TELL US HOW MUCH THE EFFECT SIZE VARIES

The forest plot of a meta-analysis typically includes a line with the summary effect size and its confidence interval, which is sometimes displayed as a diamond. Researchers sometimes assume that the confidence interval tells us how widely the effect size varies across studies. It does not. The confidence interval speaks to the precision with which we have estimated the mean effect size. It says nothing about the dispersion in effects. See Chapter 17 for a detailed discussion of this point.

### THE $I^2$ STATISTIC DOES NOT TELL US HOW MUCH THE EFFECT SIZE VARIES

Many researchers believe that the  $I^2$  index tells us how much the effect size varies, but in fact, it does not. While many readers will find this statement surprising, the proof is both simple and compelling. In this analysis,  $I^2$  was reported as 67%. Based on that, does the distribution of effects resemble panel A, B, or C? The answer is that we do not know.

There is a widespread belief that  $I^2$  values of less than 25% represent low heterogeneity; values near 50% moderate heterogeneity; and values greater than 75% high heterogeneity. Since the  $I^2$  value of 67% falls in the moderate to high interval, some researchers may expect that the dispersion in this case resembles panel B or C. That happens to be incorrect, since the dispersion in this case actually resembles panel A. However, the more important point is that given an  $I^2$  value of 67%, the heterogeneity could resemble panel A, B, C, or an infinite number of other panels.  $I^2$  does not tell us how much the effect size varies. It was never intended for that purpose and cannot provide that information except in special cases (Borenstein, 2019; Borenstein, 2020; Higgins, Hedges, & Rothstein, 2017; Huedo-Medina, Sanchez-Meca, Marin-Martinez, & Botella, 2006; Mittelbock & Heinzl, 2006; Rucker, Schwarzer, Carpenter, & Schumacher, 2008).

### WHAT $I^2$ TELLS US

If  $I^2$  does not tell us how much the effect size varies, one might ask what it does tell us. To explain that, we need to provide some background.

When we discuss a meta-analysis, we need to distinguish between *true* effects and *observed* effects. The *true* effect size in any study is the effect size that we would

observe if we could somehow enroll the entire population in the study, so that we knew the effect size with no error. By contrast, the *observed* effect size is the effect size observed in the study's sample. This serves as an estimate of the true effect size but invariably underestimates or overestimates the true effect size due to sampling error.

When we perform a meta-analysis, we work with the *observed* effect size for each study, but what we really care about is the *true* effect size for each study. As it happens, the dispersion of *observed* effects tends to exceed the dispersion of *true* effects. To understand why this is, consider what would happen if we drew ten random samples from the same population. Since all samples are estimating the same parameter (the effect size in that one population), the variance in true effects is zero by definition. Nevertheless, the variance of observed effects will be greater than zero because of sampling error. In this case,

$$V_{OBS} = V_{ERR}, \quad (19.1)$$

where  $V_{OBS}$  is the variance of observed effects, and  $V_{ERR}$  is the variance due to sampling error. The same idea applies when the variance of true effects exceeds zero. In this case, the variance of observed effects is equal to the variance of true effects plus the error variance. That is,

$$V_{OBS} = T^2 + V_{ERR}, \quad (19.2)$$

where  $T^2$  is the variance of true effects.

For the present discussion, the key point is that we have two distinct distributions. One is based on the variance of observed effects (which we see in the forest plot). The other is based on the variance of true effects (which tells us how much the effects actually vary). And, the variance of the former is greater than the variance of the latter. It would be useful to have a statistic that gives us the relationship between the two variances. That statistic is  $I^2$ , which is defined as

$$I^2 = \left( \frac{V_{TRUE}}{V_{OBS}} \right) \times 100 = \left( \frac{T^2}{V_{OBS}} \right) \times 100 = \left( \frac{V_{TRUE}}{V_{TRUE} + V_{ERR}} \right) \times 100. \quad (19.3)$$

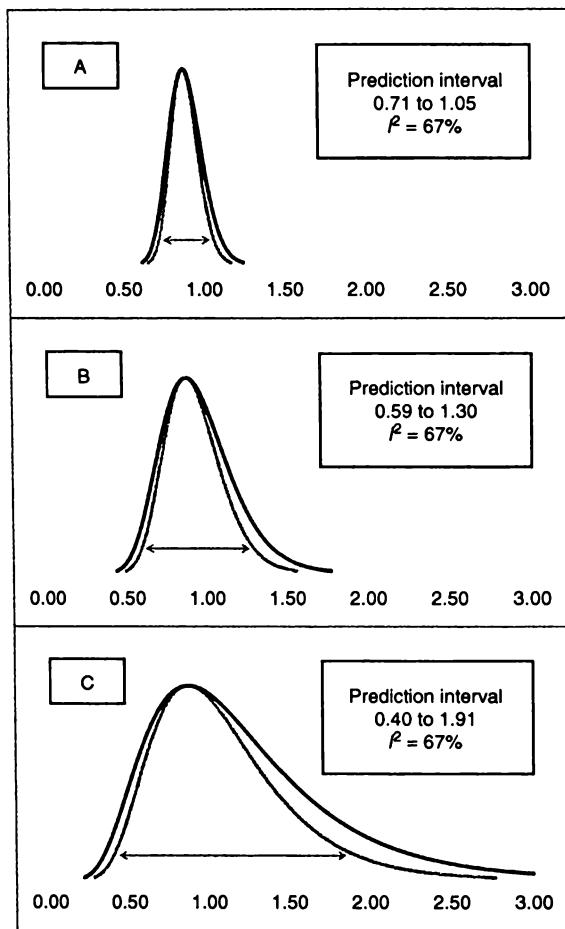
In words,  $I^2$  tells us what proportion of the observed variance is attributed to the variance in true effects rather than to sampling error (Borenstein, 2019; Borenstein, Higgins, Hedges, & Rothstein, 2017; Higgins & Thompson, 2002).

Critically,  $I^2$  is a proportion, not an absolute value. To obtain the variance of true effects ( $T^2$ ), we need to multiply  $I^2$  by the variance of observed effects. That is,

$$T^2 = I^2 \times V_{OBS}. \quad (19.4)$$

The practical implications of this formula are evident in Figure 19.2. This is based on the same example as Figure 19.1, but in this case each panel displays two distributions. The inner curve is the same curve that we saw in Figure 19.1 and reflects the dispersion of true effects, with an arrow denoting the 95% prediction interval. The outer curve represents the dispersion of observed effects.

- In panel A, the observed effects all fall within the outer curve, which is relatively narrow in this case. When we multiply this by  $I^2$ , we find that the true effects fall in the relatively narrow interval of 0.71 to 1.05, as indicated by the inner curve.



**Figure 19.2** Alcohol use and mortality. Risk ratio < 1 favors drinkers. Three possible distributions of true effects (inner) and observed effects (outer).

- In panel B, the observed effects again fall within the outer curve, which is relatively wide in this case. When we multiply this by  $I^2$ , we find that the true effects fall in the relatively wide interval of 0.59 to 1.30, as indicated by the inner curve.
- In panel C, the observed effects again fall within the outer curve, which is even wider in this case. When we multiply this by  $I^2$ , we find that the true effects fall in even wider the interval of 0.40 to 1.91, as indicated by the inner curve.

## THE $I^2$ INDEX VS. THE PREDICTION INTERVAL

If we want to know what proportion of the variance in observed effects is attributed to variance in true effects, we look at the relationship between the two curves. In all three panels, the relationship between the inner curve and the outer curve is the same, with the variance of true effects being 67% as large as the variance of observed effects. In all three cases,  $I^2$  is 67%. This is the domain of  $I^2$  – it addresses the *ratio* of true to total variance.

By contrast, if we want to know *how much* the effect size varies, we are asking for an absolute measure of dispersion. In panel A, the effects fall in the interval of 0.71 to 1.05. In panel B, they fall in the interval of 0.59 to 1.30. In panel C, they fall in the interval of 0.40 to 1.91. This is the domain of the prediction interval – it addresses the extent of the dispersion on an absolute scale (Borenstein, 2019, 2020; Borenstein *et al.*, 2017; IntHout, Ioannidis, Rovers, & Goeman, 2016).

Since  $I^2$  is a ratio, the  $I^2$  value of 67% could correspond to any of these panels. As it happens, the observed effects correspond to the outer curve in panel A, and so the true effects correspond to the inner curve in panel A. Computations are presented at the end of this chapter.

## THE PREDICTION INTERVAL

When we ask about heterogeneity in a meta-analysis, we want to know how much the effect size varies across studies. As discussed above, the  $I^2$  index does not provide this information. The index that does provide this information is the prediction interval (Borenstein, 2019; Michael Borenstein *et al.*, 2017; Chiolero, Santschi, Burnand, Platt, & Paradis, 2012; Graham & Moran, 2012; Guddat, Grouven, Bender, & Skipka, 2012; Higgins, Thompson, & Spiegelhalter, 2009; Riley, Higgins, & Deeks, 2011).

When we perform a random-effects analysis, we assume that the studies in the analysis are a random (or at least representative) sample of studies in some universe of interest, and our goal is to make inferences about that universe. The 95% prediction interval is the interval that includes the true effect size for 95% of all populations in that universe.

Figure 19.3 is a forest plot of the studies in the motivating example. The last line on the plot [A] shows the mean effect size of 0.87 with a confidence interval of 0.83 to 0.91. The confidence interval is an index of *precision*, and it speaks to the precision with which we have estimated the mean. In 95% of all analyses, the true *mean* for the universe of comparable studies will fall within the confidence interval. The confidence interval, shown here as a line, is often shown as a diamond. It has the same meaning in either case.

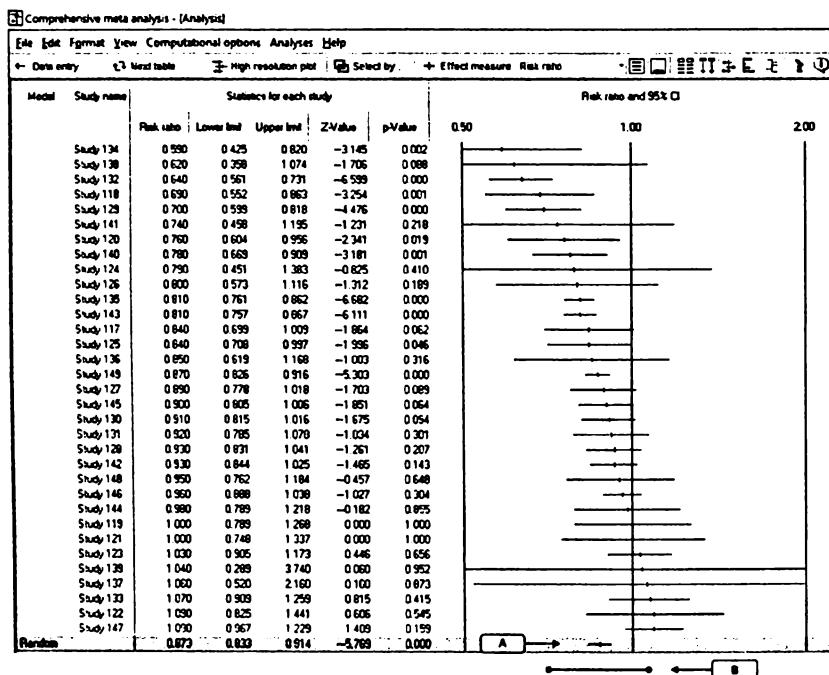


Figure 19.3 Alcohol use and mortality (Forest plot). Risk ratio < 1 favors drinkers.

The line immediately below the plot [B] shows the prediction interval of 0.71 to 1.07. The prediction interval is an index of *dispersion*, and it speaks to the heterogeneity of true effects. In some 95% of all comparable populations, the risk ratio will fall in this interval.

Researchers sometimes confuse the prediction interval with the confidence interval, but the two are entirely separate issues.

The 95% confidence interval may be estimated using

$$CI = M \pm 1.96(SE), \quad (19.5)$$

where  $SE$  is the standard *error* of the mean effect size. By contrast, the 95% prediction interval may be estimated using

$$PI = M \pm 1.96(T), \quad (19.6)$$

where  $T$  is the standard *deviation* of the effect size.

Both of these formulas are simplified versions of the formulas that we would use in practice. We use these here to highlight the difference between the confidence interval

(which is based on the standard *error*) and the prediction interval (which is based on the standard *deviation*).

In practice, one might employ the Knapp–Hartung adjustment when computing the confidence interval, as discussed in Chapter 26. Similarly, we would always recommend using the formulas discussed in Chapter 17 when computing the prediction interval. These formulas adjust the width of the interval to account for the fact that the statistics included in these formulas are estimated with error.

### PREDICTION INTERVAL IS CLEAR, CONCISE, AND RELEVANT

The prediction interval is concise and unambiguous. If we report that the risk ratio varies from 0.71 in some populations to 1.07 in others, the reader understands what this means. The prediction interval is intuitive because it reports values on the same scale as the effect size. In the motivating example, the mean effect size is a risk ratio of 0.87 and the prediction interval tells us that in most comparable populations the true risk ratio will fall between 0.71 and 1.07.

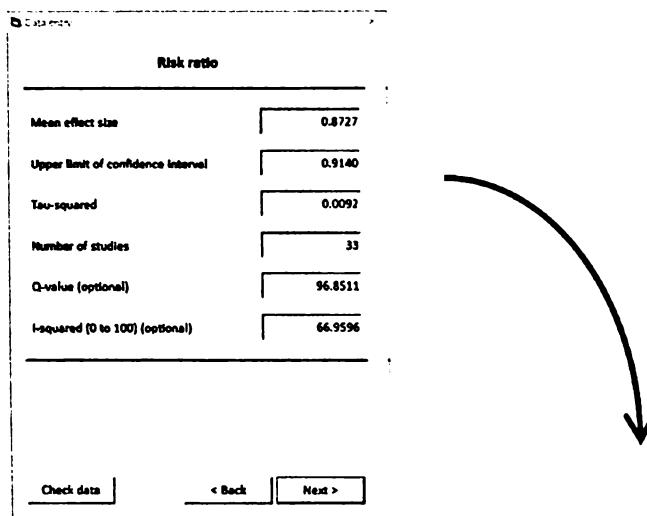
The prediction interval is on a meaningful scale. It tells us not only that the interval is a specific width, but that it ranges from one specific value to another specific value. As such, it allows us to distinguish not only between a case where the interval is 20 points wide from one where it is 40 points wide. It also allows us to distinguish between a case where those 40 points vary from trivially helpful to moderately helpful (on the one hand) vs. a case where the effects vary from harmful to helpful (on the other).

Most important, the prediction interval addresses the question that we have in mind when we ask about heterogeneity. If the analysis addresses the impact of an intervention, the prediction interval provides the information that speaks to the potential utility of that intervention.

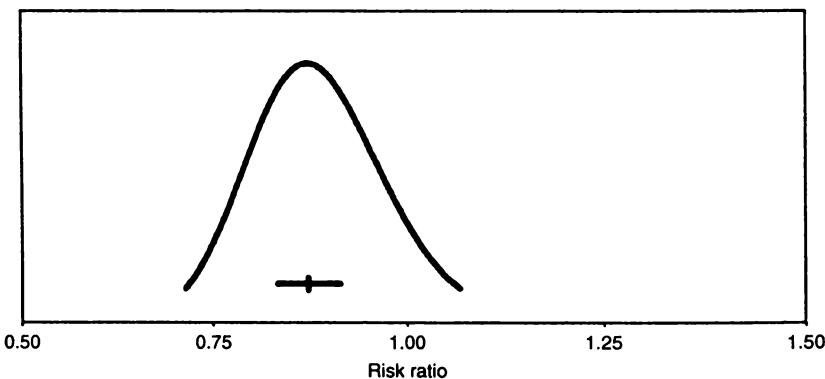
### COMPUTING THE PREDICTION INTERVAL

The formula for computing the prediction interval presented above (19.6) was intended as a conceptual formula. In practice, we need to modify this formula to account for the fact that we are working with an *estimate* of the true mean effect size and an *estimate* of the standard deviation of effects. Additionally, for some effect size indices we need to transform the estimates into log units for the computations. The relevant formulas are given in Chapter 17, and worked examples are presented in Chapter 18.

In practice, one would use a spreadsheet or computer program to compute the prediction interval. Figure 19.4 shows a program which is available on the book's website ([www.Introduction-to-Meta-Analysis.com](http://www.Introduction-to-Meta-Analysis.com)). The researcher enters the number of studies in the analysis, the risk ratio, the upper limit of the confidence interval, and tau-squared. The program generates the corresponding distribution. It also generates



### Alcohol use and all-cause mortality



The mean effect size is 0.87 with a 95% confidence interval of 0.83 to 0.91  
 The true effect size in 95% of all comparable populations falls in the interval 0.71 to 1.07

**Figure 19.4** Alcohol use and mortality (true effects). Risk ratio < 1 favors drinkers.

the caption *The true effect size in 95% of all comparable populations falls in the interval 0.71 to 1.07*, which is the prediction interval.

The formulas implemented in this program are explained in this volume, and in (Borenstein *et al.*, 2017; Higgins *et al.*, 2009; Riley *et al.*, 2011). Other approaches to computing prediction intervals are discussed in Nagashima, Noma, & Furukawa

(2019); and Wang & Lee (2019). Computational details for this example are presented at the conclusion of this section.

## HOW TO USE $I^2$

While  $I^2$  does not tell us how much the effect size varies on an absolute scale, it does provide other information, as follows.

- If  $I^2$  is zero, then all the variance in observed effects is due to sampling error. The variance in true effects is estimated as zero.
- If we are looking at a forest plot,  $I^2$  provides context for understanding that plot. If  $I^2$  is near zero, the variance of true effects is only a small fraction of that suggested by the plot. As  $I^2$  increases, that proportion increases.
- If we are working with a set of meta-analyses where the variance of observed effects is reasonably consistent, there will be a strong correlation between  $I^2$  and the absolute amount of variance. Within that context,  $I^2$  can provide information about the relative amounts of dispersion across analyses.
- The  $I^2$  statistic can be used to compare meta-analyses of the same set of data analyzed using different effect metrics. For example, raw mean differences and standardized mean differences will be associated with different amounts of heterogeneity, but it is not meaningful to compare values of  $T^2$  between the two scales. Because  $I^2$  statistic has a unit-less scale, it is legitimate to compare it between the two analyses.
- The  $I^2$  statistic is useful to statisticians who are evaluating the properties of various statistics. For example, if someone wanted to run simulations to see how statistical power is affected by the ratio of true to total variance, they could do so for various values of  $I^2$ .
- Sometimes, we do care about the proportion of variance rather than the absolute amount of variance. For example, if we have various ways of conducting studies and we want to know which have the smallest amount of sampling error,  $I^2$  is the index that allows us to address this question.

## HOW TO EXPLAIN HETEROGENEITY

Virtually all papers that report a meta-analysis include a discussion of heterogeneity which follows a standard pattern. The researchers report  $Q$ ,  $df$ , and a  $p$ -value,  $I^2$ , and  $T^2$ . None of these directly addresses the question that really matters, which is ‘What is the interval over which the effects vary?’ Ironically, the prediction interval, the one statistic that does address this question, is rarely reported.

The paragraph that follows is based on the motivating example and can be adapted for the results section of a paper. The paragraph includes all the statistics that readers (and journal editors) expect to see, but these are annotated to make it more likely that they will be interpreted correctly. Critically, the report also includes the prediction interval.

## HOW MUCH DOES THE EFFECT SIZE VARY ACROSS STUDIES?

The  $Q$ -statistic provides a test of the null hypothesis that all studies in the analysis share a common effect size. If all studies shared the same effect size, the expected value of  $Q$  would be equal to the degrees of freedom (the number of studies minus 1). The  $Q$ -value is 96.851 with 32 degrees of freedom and  $p < 0.001$ . We can reject the null hypothesis that the true effect size is the same in all these studies. The  $I^2$  statistic is 67%, which tells us that 67% of the variance in observed effects reflects variance in true effects rather than sampling error.  $T^2$ , the variance of true effect sizes, is 0.009 in log units.  $T$ , the standard deviation of true effect sizes, is 0.096 in log units. If we assume that the effects are normally distributed (in log units), we can estimate that the prediction interval for the risk ratio is 0.71 to 1.07. The true effect size for any single population will usually fall in this range.

## CAVEATS

All heterogeneity statistics will only be reliable if certain assumptions are met. In particular, we need to have a sufficient number of studies, and these studies must be a random sample of the intended universe. We also assume that the effects are normally distributed on the relevant scale. There is no consensus on what would be a sufficient number of studies to yield reliable estimates, but ten studies would be a useful minimum in most cases. With fewer than ten studies,  $T^2$  (which feeds into the prediction interval) is estimated erratically and may give rise to prediction intervals that are inappropriately narrow or unhelpfully wide. While this caveat applies to  $I^2$  and  $T^2$  as well as the prediction interval, it is of particular import for the prediction interval since researchers understand what that interval means and will actually use it in discussing the utility of an intervention.

## CONCLUSION

When we ask about heterogeneity in effects, we intend to ask how much the effect size varies across studies. We want to know the extent of the variation – *Does the effect size vary over 10 points or 50 points?* We also want to know the limits of the variation on an absolute scale – *Is the intervention always helpful, or is it helpful in some cases and harmful in others?*

The  $I^2$  index has become the primary index for reporting heterogeneity in a meta-analysis and is widely interpreted as telling us how much the effect size varies across studies. However, this interpretation is fundamentally incorrect. The  $I^2$  index is a proportion, not an absolute amount. It tells us what proportion of the variance in observed effects is attributed to variance in true effects. It does not tell us how much the effect size varies across studies.

The index that does tell us how much the effect size varies across studies is the prediction interval. This index is intuitive and concise. It reports the interval using the

same scale as the effect size itself. It gives us not only the width of the interval but the limits, so we know if the intervention is consistently helpful, or if it may be harmful in some cases. This statistic addresses the issue that we care about, and that many researchers *think* is being addressed by  $I^2$ . The inclusion of this interval in reports of heterogeneity will allow for a more informed discussion of the potential utility of any intervention and should be made common practice.

## FURTHER READING

The original papers on  $I^2$  are Higgins and Thompson (2002); Higgins, Thompson, Deeks, and Altman (2003). For a more detailed discussion of the issues raised in this section, see Borenstein *et al.* (2017).

(For related papers, see Borenstein, 2019, 2020; Coory, 2009; Higgins, 2008; Higgins *et al.*, 2009; Huedo-Medina *et al.*, 2006; IntHout *et al.*, 2016; Ioannidis, 2008a; Riley *et al.*, 2011; Rucker *et al.*, 2008).

### SUMMARY POINTS

- When we ask about heterogeneity, we want to know how the effect size varies across studies. The statistics typically reported for heterogeneity do not provide this information.
- This information is not provided in a useful format by the  $Q$ -value, nor by  $T^2$ . There is a widespread belief that the  $I^2$  index in a meta-analysis tells us how much the effect size varies across studies, but this belief is fundamentally incorrect.
- The only statistic that directly reports this information is the prediction interval. The prediction interval tells us how much the effect size varies. It tells us whether the intervention is consistently helpful, or helpful in some populations but harmful in others. This is the information that we need to make informed decisions about the potential utility of the intervention.

## THE MEANING OF $I^2$ IN FIGURE 19.2

The purpose of Figure 19.2 is to show how  $I^2$  reflects the relationship between the inner curve (true effects) and the outer curve (observed effects). In this example,  $I^2$  is 67%, which tells us that the ratio is the two variances is 0.67. However, it may not be clear how we see this in the plot. The computations for Panel A in Figure 19.2 are given in Table 19.1. In practice, we would bypass these computations and compute the prediction interval directly. This section is intended only to explain the relationship among the indices.

When we are working with risk ratios, data are converted to natural log units and all computations are performed in this metric. Therefore, most columns in the table are

**Table 19.1** Relationship between observed effects and true effects in Figure 19.2, Panel A.

	Log units			Risk ratio units	
	Mean	Variance	Standard deviation	Interval	Interval
True	-0.136154	$T^2 = 0.009222$	$T = 0.096031$	-0.324375 to 0.052067	0.722979 to 1.053446
Observed	-0.136154	$S^2 = 0.013772$	$S = 0.117356$	-0.366172 to 0.093864	0.693383 to 1.098410

in these units. After the intervals are computed, they are converted back to risk ratio units as in the right-most columns.

The  $I^2$  index is a ratio of variances,

$$I^2 = \frac{T^2}{S^2} = \frac{0.009222}{0.013772} = 0.669596 \approx 67\%. \quad (19.7)$$

To move from the variance of the outer curve to the variance of the inner curve, we would use

$$T^2 = S^2 \times I^2 = 0.013772 \times 0.669596 = 0.009222. \quad (19.8)$$

On that basis, many researchers expect that the inner curve will be 67% as wide as the outer curve. In fact, though, the ratio applies to the *variance* of the two distributions, which is a squared metric. By contrast, the distributions are in linear units, and so based on standard deviations rather than variances. Rather than work with the squared metric ( $I^2$ ), we work with the linear metric ( $I$ ).

The  $I$  index is a ratio of standard deviations,

$$I = \frac{T}{S} = \frac{0.096031}{0.117356} = 0.818288, \quad (19.9)$$

or simply

$$I = \sqrt{I^2} = \sqrt{0.669596} = 0.818288 \approx 82\%. \quad (19.10)$$

To move from the standard deviation of the outer curve to the standard deviation of the inner curve, we would use

$$T = S \times I = 0.117356 \times 0.818288 = 0.096031. \quad (19.11)$$

This is what we see in the plot – the standard deviation of the inner curve is 82% as large as the standard deviation of the outer curve.

To compute the 95% interval for observed effects we use

$$OBS_{LL} = M - 1.96(S) = -0.136154 - 1.96(0.117356) = -0.366172 \quad (19.12)$$

$$OBS_{UL} = M + 1.96(S) = 0.136154 + 1.96(0.117356) = 0.093864. \quad (19.13)$$

We then convert these log values into risk ratio units, using

$$OBS_{LL} = \exp(-0.366172) = 0.693383 \quad (19.14)$$

$$OBS_{UL} = \exp(0.093864) = 1.098410. \quad (19.15)$$

To compute the 95% interval for the inner curve at the location of the arrow, we use the same formula, but substitute the standard deviation of true effects ( $T$ ) for the standard deviation of observed effects ( $S$ ). Concretely

$$PRED_{LL} = M - 1.96(T) = -0.136154 - 1.96(0.096031) = -0.324375 \quad (19.16)$$

$$PRED_{UL} = M + 1.96(T) = -0.136154 + 1.96(0.096031) = 0.052067. \quad (19.17)$$

We then convert these log values into risk ratio units, using

$$PRED_{LL} = \exp(-0.324375) = 0.722979 \quad (19.18)$$

$$PRED_{UL} = \exp(0.052067) = 1.053446 \quad (19.19)$$

The distribution of observed effects is a hypothetical distribution that allows us to illustrate the meaning of  $I^2$ . In real life, such a distribution would only exist if the error variance was identical for all studies, which is never the case.

The formulas used here to compute the prediction interval are only intended for the purpose of illustrating the relationship between the curves. For that purpose, we wanted to use a formula that isolates the difference between the distribution of true effects vs. observed effects. By contrast, to compute the prediction interval in practice, we would use the formulas presented earlier that use the  $t$  distribution rather than the  $Z$  distribution, and that take into account the error variance in estimating the mean. In this example, the prediction interval based on the correct formulas (0.71 to 1.07) is only slightly wider than the one based on the naïve formulas (0.72 to 1.05). This is true in this example because we have a substantial number of studies and a precise estimate of the mean. However, it would be a mistake to generalize from this example and assume that we can always use the naïve formulas. Often, the difference between the naïve formula and the correct formula will be substantial, and so we should always use the latter.



# Classifying Heterogeneity as Low, Moderate, or High

---

### Introduction

Interest should generally focus on an index of absolute heterogeneity  
The classifications lead themselves to mistakes of interpretation  
Classifications focus attention in the wrong direction

---

### INTRODUCTION

In recent years, researchers have developed the practice of classifying heterogeneity as being low, moderate, or high based on the value of  $I^2$ . For example, some classify an  $I^2$  value of 25% or less as low; 50% as moderate; and 75% or more as high. These classifications were proposed by one of the authors of this book, with specific reference to values that one might expect to see in meta-analyses of clinical trials in the Cochrane Database of Systematic Reviews. Unfortunately, such interpretations are widely applied inappropriately. Here we argue that the idea of classifying heterogeneity based on  $I^2$  should be strongly discouraged.

### INTEREST SHOULD GENERALLY FOCUS ON AN INDEX OF ABSOLUTE HETEROGENEITY

The first reason that we should not use this classification system is that when we talk about heterogeneity as being low, moderate, or high, we are talking about the *absolute* amount of heterogeneity. It follows that any classification scheme intended to address this goal must be based on an index that reports the *absolute* amount of heterogeneity. It could not be based on  $I^2$  since this index does not tell us how much the effects vary. Consider the following two examples.

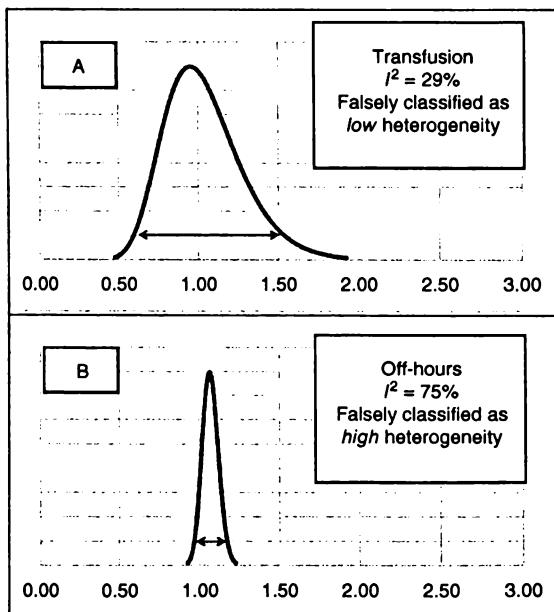
Holst, Petersen, Haase, Perner, and Wetterslev (2015) looked at the relationship between two treatment conditions (restrictive vs. liberal criterion for blood transfusion) and outcome. An odds ratio less than 1.0 would indicate that patients treated with

the restrictive strategy were more likely to have a good outcome. The odds ratio was 0.96 with a confidence interval of 0.78 to 1.18. We will refer to this as the transfusion analysis.

Sorita *et al.* (2014) looked at the relationship between two treatment conditions (patients who presented at the hospital with acute MI during normal hours vs. off hours) and short-term mortality. The mean effect size is an odds ratio of 1.06 with a confidence interval of 1.04 to 1.09, indicating that patients who presented during off hours had a higher risk of death. We will refer to this as the off-hours analysis.

In the transfusion analysis, the  $I^2$  statistic was 29%. In the off-hours analysis, the  $I^2$  statistic was 75%. On that basis, many would assume that the effects varied more widely in the second analysis as compared with the first. As it happens, the opposite is true. In Figure 20.1, the top panel shows the distribution of effects for the transfusion analysis, and the bottom panel shows the distribution of effects for the off-hours analysis. When  $I^2$  was 29%, the effects vary from 0.60 to 1.51. By contrast, when  $I^2$  was 75%, the effects vary from 0.97 to 1.17. Thus, the *lower* value of  $I^2$  corresponds to the *greater* amount of dispersion.

Additionally, if we were to use these values of  $I^2$  to classify the dispersion, the dispersion in the top panel would be classified as low, while the dispersion in the bottom panel would be classified as high. This is clearly misleading since the dispersion in the top panel (classified as low) is roughly five times as wide as that in the bottom panel (classified as high).



**Figure 20.1** True effects for two meta-analyses.

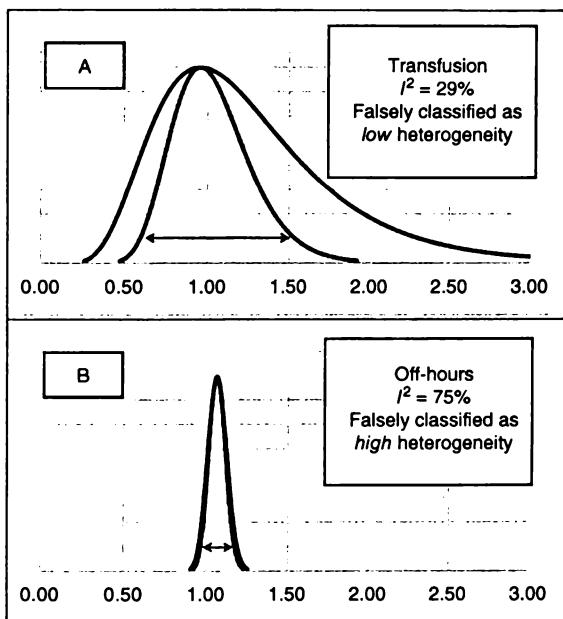


Figure 20.2 True effects (inner) and observed effects (outer) for two meta-analyses.

This is inexplicable to persons who believe that  $I^2$  is an index of absolute dispersion. However, it makes sense to those who understand that  $I^2$  tells us the relationship between the true effects and the observed effects. This is illustrated in Figure 20.2. In the top panel, the variance of true effects is only 29% as large as the variance of observed effects, and so  $I^2$  is 29%. In the bottom panel, the variance of true effects is 75% as large as the variance of observed effects, and so  $I^2$  is 75%. As it happens, the *observed* variance in the top panel is substantially greater than the *observed* variance in the bottom panel. In this case, 29% of the (relatively little) observed variance in the top panel turns out to be more than 75% of the (relatively large) observed variance in the bottom panel.

The use of the misleading classifications has serious implications. In the transfusion analysis, based on the (incorrect) classification of low heterogeneity one might assume that the mean effect size applies to all populations. In fact, though, the plot suggests that there will be some populations where the restrictive approach is much more effective, and other populations where the conservative approach is much more effective. Conversely, in the off-hours analysis, based on the (incorrect) classification of high heterogeneity, there would be some populations where the additional risk was trivial and some where it was substantial. In fact, though, the effect was relatively consistent across all comparable populations.

To be clear, these examples were not selected at random. If we were to choose two analyses at random, the analysis with the higher value of  $I^2$  would have more variance

*on average.* Rather, these examples were selected to make the point that we cannot classify heterogeneity as being low, moderate, or high based on  $I^2$ . The  $I^2$  index does not tell us how much the effects vary. Indeed, as in this example, it cannot reliably tell us which of two analyses has more dispersion.

Note that the distributions of effects are assumed to be symmetric in log units. They appear to be skewed because the plots use the risk ratio rather than the log risk ratio on the X-axis.

### THE CLASSIFICATIONS LEAD THEMSELVES TO MISTAKES OF INTERPRETATION

While it would be possible to develop objective standards for what constitutes a low, moderate, or high amount of heterogeneity from a statistical perspective, the fact remains that researchers and readers understand these terms in the colloquial sense. If we report that the heterogeneity is ‘low’, this will be taken to mean that the impact of the intervention is consistent in a clinical or substantive sense. However, the clinical interpretation depends on the context. Variation in effects that we might consider to be trivial in one context might be clinically important in another. And even for a given context, readers will have different ideas of what ‘low’ or ‘moderate’ means.

By contrast, the prediction interval reports the actual range of effects and ensures that all readers have a common understanding of the dispersion.

### CLASSIFICATIONS FOCUS ATTENTION IN THE WRONG DIRECTION

Finally, the use of a classification system encourages researchers to focus primarily on the amount of dispersion, rather than the clinical or substantive implications of the dispersion. When considering the potential utility of an intervention, we need to consider the larger picture, which requires a synthesis of the mean effect size and the heterogeneity. For example, rather than ask about the amount of variation, we might ask if the effect is consistently helpful, or if it is helpful in some cases and harmful in others.

We return to that in chapter 24

#### SUMMARY POINTS

- The practice of classifying heterogeneity as being low, moderate, or high based on the value of  $I^2$  should be emphatically discouraged.
- The  $I^2$  statistic does not tell us how much the effect size varies, and therefore, a classification based on  $I^2$  cannot tell us how much the effect size varies.
- Dispersion classified as ‘low’ in one meta-analysis could be substantially greater than dispersion classified as ‘high’ in another meta-analysis.

# Explaining Heterogeneity



# Subgroup Analyses

---

Introduction

Fixed-effect model within subgroups

Computational models

Random effects with separate estimates of  $\tau^2$

Random effects with pooled estimate of  $\tau^2$

The proportion of variance explained

Mixed-effects model

Obtaining an overall effect in the presence of subgroups

---

## INTRODUCTION

In the preceding chapters we explained how to quantify heterogeneity in the effect size across studies. When the effect size varies across studies, we generally want to understand what variables are associated with that variation.

In this chapter we show how meta-analysis can be used to compare the mean effect for different subgroups of studies (akin to analysis of variance in a primary study). In the next chapter we show how meta-analysis can be used to assess the relationship between study-level covariates and effect size (akin to multiple regression in primary studies). In chapter 23 we provide additional context for both of these methods.

Consider the following examples.

- We anticipate that a class of drugs reduces the risk of death in patients with cardiac arrhythmia, but we hypothesize that the magnitude of the effect depends on whether the condition is acute or chronic. We want to determine whether the drug is effective for each kind of patient, and also to determine whether the effect differs in the two.
- Our meta-analysis includes ten studies that used proper randomization techniques and ten that did not. Before computing a summary effect across all 20 studies we want to compute the effect for each group of ten, and determine if the effect size is related to the kind of randomization employed in the study.

- We anticipate that forest management reduces the destruction of tree stands by insect pests, but we hypothesize that the magnitude of the effect depends on the diversity of trees in the stand. We want to determine whether forest management is effective in reducing destruction for both single species and mixed stands, and also to determine whether the effect differs in the two.
- We have data from ten studies that looked at the impact of tutoring on math scores of ninth-grade students. Five of the studies used one variant of the intervention while five used another variant. We anticipate that both variants are effective, and our primary goal in the analysis is to determine whether one is *more effective* than the other.

We shall pursue the last of these examples (the impact of tutoring on math) throughout this chapter. The effect size in this example is the standardized mean difference between groups (Hedges'  $g$ ) but the same formulas would apply for any effect size index. As always, if we were working with odds ratios or risk ratios all values would be in log units, and if we were working with correlations all values would be in Fisher's  $z$  units.

Assume all the studies used the same design, with some students assigned to be tutored and others to a control condition. In some studies (here called  $A$ ) students were tutored once a week while in the others ( $B$ ) students were tutored twice a week. Our goal is to compare the impact of the two protocols to see if either intervention is more effective than the other.

Note. In this example we will be comparing the effect in one subgroup of studies versus the effect in a second subgroup of studies. The ideal scenario would be to have studies that directly compare the two variants of the intervention, since this would remove the potential for confounds and also reduce the error term. We assume that such studies are not available to us.

### How this chapter is organized

We present three computational *models*. These are (a) fixed-effect, (b) random-effects using separate estimates of  $\tau^2$ , and (c) random-effects using a pooled estimate of  $\tau^2$ .

For each of the three models we present three *methods* for comparing the subgroups. These are (1) the  $Z$ -test, (2) a  $Q$ -test based on analysis of variance, and (3) a  $Q$ -test for heterogeneity.

The three statistical *models*, crossed with the three computational *methods*, yield a total of nine possible combinations. These are shown in Box 21.1, which serves as a roadmap for this chapter. Readers who want to get a sense of the issues quickly may find it easier to read the introduction and method 1 for each model, and return later to methods 2 and 3.

The dataset and all computations are available on the book's website.

**BOX 21.1 ROADMAP**

	Introduction	Method 1	Method 2	Method 3
Model		Z-test	$Q$ -test based on ANOVA	$Q$ -test for heterogeneity
Fixed-effect	Page 163	Page 167	Page 168	Page 170
Random-effects with separate estimates of $\tau^2$	Page 174	Page 178	Page 179	Page 180
Random-effects with pooled estimate of $\tau^2$	Page 181	Page 185	Page 186	Page 187

**FIXED-EFFECT MODEL WITHIN SUBGROUPS**

A forest plot of the Tutoring studies is shown in Figure 21.1. The five *A* studies (at the top) have effect sizes (Hedges'  $g$ ) in the approximate range of 0.10 to 0.50. The five *B* studies (below) have effect sizes in the approximate range of 0.45 to 0.75.

The combined effect for the *A* studies (represented by the first diamond) is 0.32 with a 95% confidence interval of plus/minus 0.11. The combined effect for the *B* studies (represented by the second diamond) is 0.61 with a 95% confidence interval of plus/minus 0.12. Our goal, then, is to compare these two effects.

If we were working with a primary study (with Thornhill, Kendall, etc. being persons in treatment *A*, and Jeffries, Fremont, etc. being persons in treatment *B*), we would compute the mean and variance for each treatment, and our options for comparing these means would be clear. For example, we could perform a *t*-test to assess the difference between means relative to the standard error of the difference. Or, we could use analysis of variance to assess the variance among groups means relative to the variance within groups.

In meta-analysis we are working with subgroups of *studies* rather than groups of *subjects*, but will follow essentially the same approach, using a variant of the *t*-test or a variant of analysis of variance to compare the subgroup means. For this purpose we need to perform two tasks.

- Compute the mean effect and variance for each subgroup.
- Compare the mean effect across subgroups.

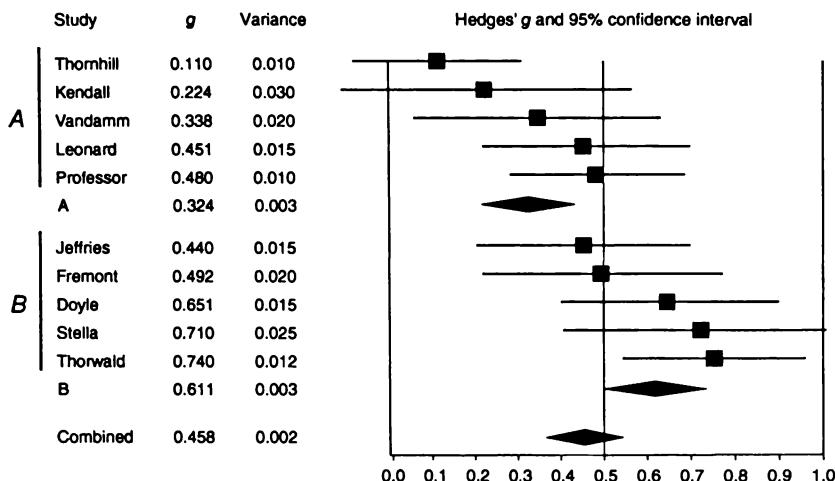


Figure 21.1 Fixed-effect model – studies and subgroup effects.

### Computing the summary effects

In Table 21.1 the data for the *A* studies are displayed at the top, and data for the *B* studies are displayed toward the bottom.

To compute the summary effects we use the same formulas that we introduced for a single group (11.2) to (11.10). The summary effect for subgroup *A* is computed

Table 21.1 Fixed effect model – computations.

Study	Effect size <i>Y</i>	Variance within <i>V<sub>Y</sub></i>	Variance between <i>T<sup>2</sup></i>	Variance total <i>V</i>	Weight <i>W</i>	Calculated quantities			
						<i>WY</i>	<i>WY<sup>2</sup></i>	<i>W<sup>2</sup></i>	
<i>A</i>	Thornhill	0.110	0.0100	0.0000	0.0100	100.000	11.000	1.210	10000.000
	Kendall	0.224	0.0300	0.0000	0.0300	33.333	7.467	1.673	1111.111
	Vandamm	0.338	0.0200	0.0000	0.0200	50.000	16.900	5.712	2500.000
	Leonard	0.451	0.0150	0.0000	0.0150	66.667	30.067	13.560	4444.444
	Professor	0.480	0.0100	0.0000	0.0100	100.000	48.000	23.040	10000.000
Sum A						350.000	113.433	45.195	28055.556
<i>B</i>	Jeffries	0.440	0.0150	0.0000	0.0150	66.667	29.333	12.907	4444.444
	Fremont	0.492	0.0200	0.0000	0.0200	50.000	24.600	12.103	2500.000
	Doyle	0.651	0.0150	0.0000	0.0150	66.667	43.400	28.253	4444.444
	Stella	0.710	0.0250	0.0000	0.0250	40.000	28.400	20.164	1600.000
	Thorwald	0.740	0.0120	0.0000	0.0120	83.333	61.667	45.633	6944.444
Sum B						306.667	187.400	119.061	19933.333
Sum						656.667	300.833	164.255	47988.889

using values from the row marked *Sum A*. The summary effect for subgroup *B* is computed using values from the row marked *Sum B*. The summary effect for all studies is computed using values from the row marked *Sum*.

**Computations (fixed effect) for the *A* studies**

$$M_A = \frac{113.433}{350.000} = 0.3241,$$

$$V_{M_A} = \frac{1}{350.000} = 0.0029,$$

$$SE_{M_A} = \sqrt{0.0029} = 0.0535,$$

$$LL_{M_A} = 0.3241 - 1.96 \times 0.0535 = 0.2193,$$

$$UL_{M_A} = 0.3241 + 1.96 \times 0.0535 = 0.4289,$$

$$Z_A = \frac{0.3241}{0.0535} = 6.0633,$$

$$p(Z_A) < 0.0001,$$

$$Q_A = 45.195 - \left( \frac{113.433^2}{350.000} \right) = 8.4316 \quad (21.1)$$

$$p(Q = 8.4316, df = 4) = 0.0770,$$

$$C_A = 350.000 - \frac{28055.556}{350.000} = 269.8413,$$

$$T_A^2 = \frac{8.4316 - 4}{269.8413} = 0.0164,$$

and

$$I_A^2 = \left( \frac{8.4316 - 4}{8.4316} \right) \times 100 = 52.5594.$$

**Computations (fixed effect) for the *B* studies**

$$M_B = \frac{187.400}{306.667} = 0.6111,$$

$$V_{M_B} = \frac{1}{306.667} = 0.0033,$$

$$SE_{M_B} = \sqrt{0.0033} = 0.0571,$$

$$LL_{M_B} = 0.6111 - 1.96 \times 0.0571 = 0.4992,$$

$$UL_{M_B} = 0.6111 + 1.96 \times 0.0571 = 0.7230,$$

$$Z_B = \frac{0.6111}{0.0571} = 10.7013,$$

$$p(Z_B) < 0.0001,$$

$$Q_B = 119.011 - \left( \frac{187.400^2}{306.667} \right) = 4.5429, \quad (21.2)$$

$$p(Q = 4.5429, df = 4) = 0.3375,$$

$$C_B = 306.667 - \frac{19933.333}{306.667} = 241.667,$$

$$T_B^2 = \frac{4.5429 - 4}{241.667} = 0.0022,$$

and

$$I_B^2 = \left( \frac{4.5429 - 4}{4.5429} \right) \times 100 = 11.9506.$$

#### *Computations (fixed effect) for all ten studies*

$$M = \frac{300.833}{656.667} = 0.4581. \quad (21.3)$$

$$V_M = \frac{1}{656.667} = 0.0015, \quad (21.4)$$

$$SE_M = \sqrt{0.0015} = 0.0390,$$

$$LL_M = 0.4581 - 1.96 \times 0.0390 = 0.3816,$$

$$UL_M = 0.4581 + 1.96 \times 0.0390 = 0.5346,$$

$$Z = \frac{0.4581}{0.0390} = 11.7396,$$

$$p(Z) < 0.0001,$$

$$Q = 164.255 - \left( \frac{300.833^2}{656.667} \right) = 26.4371,$$

$$p(Q = 26.4371, df = 9) = 0.0017,$$

$$C = 656.667 - \frac{47988.889}{656.667} = 583.5871,$$

$$T^2 = \frac{26.4371 - 9}{583.5871} = 0.0299, \quad (21.5)$$

and

$$I^2 = \left( \frac{26.4371 - 9}{26.4371} \right) \times 100 = 65.96.$$

The statistics computed above are summarized in Table 21.2.

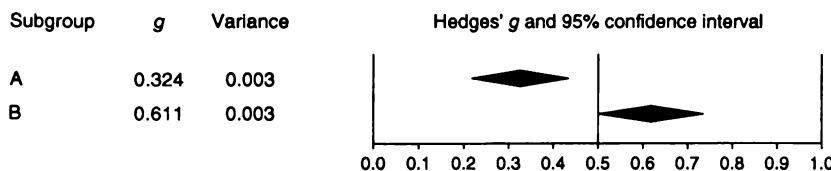
#### **Comparing the effects**

If we return to Figure 21.1 and excerpt the diamonds for the two subgroups we get Figure 21.2. The mean effect size for subgroups A and B are 0.324 and 0.611, with variances of 0.003 and 0.003.

Our goal is to compare these two mean effects, and we describe three ways that we can proceed. These approaches are algebraically equivalent, and (it follows) yield the

**Table 21.2** Fixed-effect model – summary statistics.

	A	B	Combined
Y	0.3241	0.6111	0.4581
V	0.0029	0.0033	0.0015
SEY	0.0535	0.0571	0.0390
LLY	0.2193	0.4992	0.3816
ULY	0.4289	0.7230	0.5346
Z	6.0633	10.7013	11.7396
p2	0.0000	0.0000	0.0000
Q	8.4316	4.5429	26.4371
df	4.0000	4.0000	9.0000
p-value	0.0770	0.3375	0.0017
Numerator	4.4316	0.5429	17.4371
C	269.8413	241.6667	583.5871
T <sup>2</sup>	0.0164	0.0022	0.0299
I <sup>2</sup>	52.5594	11.9506	65.9569

**Figure 21.2** Fixed-effect – subgroup effects.

same  $p$ -value. Our goal in presenting three approaches is to provide insight into the process.

### Comparing **A** versus **B**: a Z-test (Method 1)

Since there are only two subgroups here, we can work directly with the mean difference in effect sizes. In a primary study, if we wanted to compare the means in two groups we would perform a  $t$ -test. In meta-analysis the mean and variance are based on *studies* rather than *subjects* but the logic of the test is the same.

Concretely, let  $\theta_A$  and  $\theta_B$  be the true effects underlying groups *A* and *B*, let  $M_A$  and  $M_B$  be the estimated effects, and let  $V_{MA}$  and  $V_{MB}$  be their variances. If we use  $Diff$  to refer to the difference between the two effects, and elect to subtract the mean of *A* from the mean of *B*,

$$Diff = M_B - M_A,$$

the test statistic to compare the two effects is

$$Z_{Diff} = \frac{Diff}{SE_{Diff}}, \quad (21.6)$$

where

$$SE_{Diff} = \sqrt{V_{M_A} + V_{M_B}}. \quad (21.7)$$

Under the null hypothesis that the true effect size  $\theta$  is the same for both groups,

$$H_0 : \theta_A = \theta_B, \quad (21.8)$$

$Z_{Diff}$  would follow the normal distribution. For a two-tailed test, the  $p$ -value is given by

$$p = 2[1 - (\Phi(|Z|))], \quad (21.9)$$

where  $\Phi(Z)$  is the standard normal cumulative distribution.

In the running example,

$$Diff = 0.6111 - 0.3241 = 0.2870,$$

$$SE_{Diff} = \sqrt{0.0029 + 0.0033} = 0.0782,$$

$$Z_{Diff} = \frac{0.2870}{0.0782} = 3.6691,$$

and

$$p = 2[1 - (\Phi(|3.6691|))] = 0.0002.$$

The two-tailed  $p$ -value corresponding to  $Z_{Diff} = 3.6691$  is 0.0002. This tells us that the treatment effect is probably not the same for the  $A$  studies as for the  $B$  studies. In Excel, the function to compute a 2-tailed  $p$ -value for  $Z$  is  $5(1-(NORMSDIST(ABS(Z))))*2$ . Here,  $=(1-(NORMSDIST(ABS(3.6691))))*2$  will return the value 0.0002.

### Comparing $A$ with $B$ : a $Q$ -test based on analysis of variance (Method 2)

In a primary study, the  $t$ -test can be used to compare the means in *two* groups, but to compare means in more than two groups we use analysis of variance. Concretely, we partition the total variance (of all subjects about the grand mean) into variance within groups (of subjects about the means of their respective groups) and variance between groups (of group means about the grand mean). We then test these various components of variance for statistical significance, with the last (variance between groups) addressing the hypothesis that effect size differs as function of group membership.

In meta-analysis the means are based on *studies* rather than *subjects* but the logic of the test is the same. Specifically, we compute the following quantities (where  $SS$  is the sum of squared deviations).

- $Q_A$ , the weighted  $SS$  of all  $A$  studies about the mean of  $A$ .
- $Q_B$ , the weighted  $SS$  of all  $B$  studies about the mean of  $B$ .
- $Q_{within}$ , the sum of  $Q_A$  and  $Q_B$ .
- $Q_{bet}$ , the weighted  $SS$  of the subgroup means about the grand mean.
- $Q$ , the weighted  $SS$  of all effects about the grand mean.

We may write  $Q_{within} = Q_A + Q_B$ , to represent the sum of within-group weighted SS, or more generally, for  $p$  subgroups,

$$Q_{within} = \sum_{j=1}^p Q_j. \quad (21.10)$$

In the running example

$$Q_{within} = 8.4316 + 4.5429 = 12.9745. \quad (21.11)$$

The weighted SS are additive, such that  $Q = Q_{within} + Q_{bet}$ . Therefore,  $Q_{bet}$  can be computed as

$$Q_{bet} = Q - Q_{within}. \quad (21.12)$$

Under the null hypothesis that the effect size  $\theta$  is the same for all groups, 1 to  $p$ ,  $Q_{bet}$  would be distributed as chi-squared with degrees of freedom equal to  $p - 1$ .

In the running example,

$$Q_{bet} = 26.4371 - 12.9745 = 13.4626 \quad (21.13)$$

Each  $Q$  statistic is evaluated with respect to the corresponding degrees of freedom. In the running example (Table 21.3),

- The ‘Total’ line tells us that for the full group of ten studies the variance is statistically significant ( $Q = 26.4371$ ,  $df = 9$ ,  $p = 0.0017$ ).
- The ‘Within’ line tells us that the variance within groups (averaged across groups) is not statistically significant ( $Q_{within} = 12.9745$ ,  $df = 8$ ,  $p = 0.1127$ ).
- The ‘Between’ line tells us that the difference between groups (the combined effect for  $A$  versus  $B$ ) is statistically significant ( $Q_{bet} = 13.4626$ ,  $df = 1$ ,  $p = 0.0002$ ), which means that the effect size is related to the frequency of tutoring.
- At a finer level of detail, neither the variance within subgroup  $A$  ( $Q_A = 8.4316$ ,  $df = 4$ ,  $p = 0.0770$ ) nor within subgroup  $B$  ( $Q_B = 4.5429$ ,  $df = 4$ ,  $p = 0.3375$ ) is statistically significant.

As always, the absence of statistical significance (here, within subgroups) means only that we cannot rule out the hypothesis that the studies share a common effect size, and it does not mean that this hypothesis has been proven.

In Excel, the function to compute a  $p$ -value for  $Q$  is =CHIDIST( $Q$ ,  $df$ ). For the test of  $A$  versus  $B$ , =CHIDIST(13.4626, 1) returns 0.0002.

**Table 21.3** Fixed-effect model – ANOVA table.

	$Q$	$df$	$p$	Formula
A	8.4316	4	0.0770	19.1
B	4.5429	4	0.3375	19.2
Within	12.9745	8	0.1127	19.11
Between	13.4626	1	0.0002	19.13
Total	26.4371	9	0.0017	19.5

### Comparing A versus B: a Q-test for heterogeneity (Method 3)

The test we just described can be derived in a different way. We can think of the effect sizes for subgroups *A* and *B* as single studies (if we extract the two subgroup lines and the total line from Figure 21.1 and replace the diamonds with squares, to represent these as if they were studies, we get Figure 21.3). Then, we can test these ‘studies’ for heterogeneity, using precisely the same formulas that we introduced earlier (Chapter 16) to test the dispersion of single studies about the summary effect.

Concretely, we start with two ‘studies’ with effect sizes of 0.324 and 0.611, and variance of 0.003 and 0.003. Then, we apply the usual meta-analysis methods to compute *Q* (see Table 21.4).

In this example,

$$M = \frac{300.833}{656.667} = 0.4581, \quad (21.14)$$

$$V_M = \frac{1}{656.667} = 0.0015, \quad (21.15)$$

$$Q = 151.281 - \left( \frac{300.833^2}{656.667} \right) = 13.4626,$$

$$df = 2 - 1 = 1,$$

and

$$p(Q = 13.4626, df = 1) = 0.0002.$$

where *Q* represents the weighted sum of squares for studies *A* and *B* about the grand mean. For *Q* = 13.4626 and *df* = 1, the *p*-value is 0.0002.

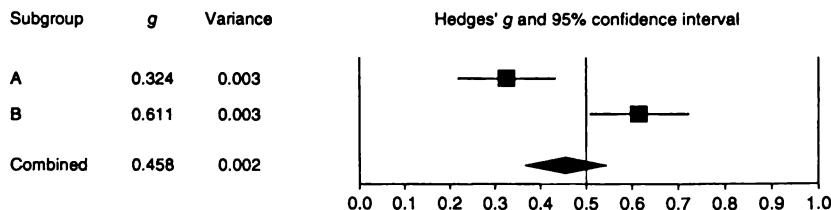


Figure 21.3 Fixed-effect model – treating subgroups as studies.

Table 21.4 Fixed-effect model – subgroups as studies.

Study	Effect size <i>Y</i>	Variance within <i>V<sub>y</sub></i>	Variance between <i>T<sup>2</sup></i>	Variance total <i>V</i>	Weight <i>W</i>	Calculated quantities		
						<i>WY</i>	<i>WY<sup>2</sup></i>	<i>W<sup>2</sup></i>
A	0.3241	0.0029	0.0000	0.0029	350.000	113.433	36.763	122500.000
B	0.6111	0.0033	0.0000	0.0033	306.667	187.400	114.518	94044.444
					656.667	300.833	151.281	216544.444

In Excel, the function to compute a  $p$ -value for  $Q$  is  $=\text{CHIDIST}(Q, df)$ , and  $=\text{CHIDIST}(13.4626, 1)$  returns 0.0002.

## Summary

We presented three methods for comparing the effect size across subgroups. One method was to use a  $Z$ -test to compare the two effect sizes directly. Another was to use a  $Q$ -test to partition the variance, and test the between-subgroups portion of the variance. A third was to use a  $Q$ -test to assess the dispersion of the summary effects about the combined effect. All the methods assess the difference in subgroup effects relative to the precision of the difference (or the variance across subgroups effects relative to the variance within subgroups).

As noted earlier, the methods are mathematically equivalent. The two methods that report  $Q$ , report the same value for  $Q$  (13.4626). When there is one degree of freedom (so that we can use either a  $Z$ -test or a  $Q$ -test)  $Z$  is equal to the square root of  $Q$ . In our example, the method that reports  $Z$ , reports a value of  $Z = 3.6691$ , which is equal to the square root of  $Q$ . All three methods yield a  $p$ -value of 0.0002.

## Quantify the magnitude of the difference

The  $Z$ -test and the  $Q$ -tests address the question of *statistical*, rather than *clinical* significance. In addition to reporting the test of significance, one should generally report an estimate of the effect size, which in this context is the difference in mean effect between the two subgroups. For subgroups  $A$  and  $B$ , if we elect to subtract the mean of  $A$  from the mean of  $B$ , the difference is

$$\text{Diff} = M_B - M_A. \quad (21.16)$$

The 95% confidence interval is estimated by

$$LL_{\text{Diff}} = \text{Diff} - 1.96 \times SE_{\text{Diff}}, \quad (21.17)$$

and

$$UL_{\text{Diff}} = \text{Diff} + 1.96 \times SE_{\text{Diff}}, \quad (21.18)$$

where the standard error was defined in (21.7). If we had more than two subgroups, we could repeat this procedure for all pairs of subgroups. In the running example the difference in effects (which we have defined as  $B$  minus  $A$ ) and its 95% confidence interval are estimated as

$$\text{Diff} = 0.6111 - 0.3241 = 0.2870,$$

$$SE_{\text{Diff}} = \sqrt{0.0029 + 0.0033} = 0.0782,$$

$$LL_{\text{Diff}} = 0.2870 - 1.96 \times 0.0782 = 0.1337,$$

and

$$UL_{\text{Diff}} = 0.2870 + 1.96 \times 0.0782 = 0.4403.$$

In words, the true difference between the effect in the subgroup *A* studies, as opposed to the subgroup *B* studies, probably falls in the range of 0.13 to 0.44.

## COMPUTATIONAL MODELS

In Part 3 of this volume we discussed the difference between a fixed-effect model and a random-effects model. Under the fixed-effect model we assume that the true effect is the same in all studies. By contrast, under the random-effects model we allow that the true effect may vary from one study to the next. This difference has implications for the way that weights are assigned to the studies, which affects both the summary effect and its standard error.

When we introduced these two models we were working with a single set of studies. Now, we are working with more than one subgroup of studies (in the running example, *A* and *B*) but the same issues apply. Under the fixed-effect model we assume that all studies in subgroup *A* share a common effect size and that all studies in subgroup *B* share a common effect size. By contrast, under the random-effects model we allow that there may be some true variation of effects within the *A* studies and within the *B* studies.

When we initially discussed the fixed-effect model we used the example of a pharmaceutical company that enrolled 1000 patients for a clinical trial and divided them among ten cohorts of 100 patients each (page 76). These ten cohorts were known to be identical in all important respects, and so it was reasonable to assume that the true effect would be the same for all ten studies. When we presented this example we noted that the conditions described (of all the studies being performed by the same researchers using the same population and methods) are rare in systematic reviews, and that in most cases the random-effects model will be more plausible than the fixed-effect.

We can expand the pharmaceutical example to apply to subgroups if we assume that five of the studies will compare *Drug A* versus placebo, and the other five will compare *Drug B* versus placebo. Within the five *Drug A* studies and within the five *Drug B* studies there should be a single true effect size, and so in this case it would be correct to use the fixed-effect model within subgroups. However, the same caveat applies here, in that this kind of systematic review, where all studies are performed by the same researchers using the same population and methods, is very rare. In the vast majority of systematic reviews these conditions will not hold, and a random-effects analysis would be a better fit for the data.

For example, in the tutoring analysis it seems plausible that the distinction between the two interventions (one hour versus two hours a week) captures some, *but not all*, of the true variation among effects. Within either subgroup of studies (*A* or *B*) there are probably differences from study to study in the motivation of the students, or the dedication of the teachers, the details of the protocol, or other factors, such that the true effect differs from study to study. If these differences do exist, and can have an impact on the effect size, then the random-effects model is a better match than the fixed-effect.

When we use the random-effects model, the impact on the summary effect within subgroups will be the same as it had been when we were working with a single population. The weights assigned to each study will be more moderate than they had

been under the fixed-effect model (large studies will lose impact while small studies gain impact). And, the variance of the combined effect will increase.

## **$T^2$ should be computed within subgroups**

To apply the random-effects model we need to estimate the value of  $\tau^2$ , the variance of true effect sizes across studies. Since  $\tau^2$  is defined as the true variance in effect size among a set of studies, its value will differ depending on how we define the set.

If we were to define the set as all studies irrespective of which subgroup they belong to, with  $\tau^2$  based on the dispersion of all studies from the grand mean,  $\tau^2$  would tend to be relatively large. By contrast, if we define the set as all studies *within* a subgroup, with  $\tau^2$  based on the dispersion of the *A* studies from the mean of *A* and of the *B* studies from the mean of *B*,  $\tau^2$  would tend to be relatively small (especially if the *A* studies and the *B* studies do represent distinct clusters, as we have hypothesized).

Since our goal is to estimate the mean and sampling distribution of subgroup *A*, and to do the same for subgroup *B*, it is clearly the variance *within* subgroups that is relevant in the present context. Put simply, if some of the variance in effect sizes can be explained by the type of intervention, then this variance is not a factor in the sampling distribution of studies within a subgroup (where only one intervention was used). Therefore, we always estimate  $\tau^2$  within subgroups.

### **To pool or not to pool**

When we estimate  $\tau^2$  within subgroups of studies, the estimate is likely to differ from one subgroup to the next. In the running example, the estimate of  $\tau^2$  in subgroup *A* was 0.016, while in subgroup *B* it was 0.002. We have the option to pool the within-group estimates of  $\tau^2$  and apply this common estimate to all studies. Alternatively, we can apply each subgroup's estimate of  $\tau^2$  to the studies in that subgroup.

Note. As a shorthand we refer to pooling the estimates of  $\tau^2$ . In fact, though, what we actually pool are  $Q$ ,  $df$ , and  $C$ , and then estimate  $\tau^2$  from these pooled values (see (21.38)).

The decision to pool (or not) depends on the following. If we assume that the true study-to-study dispersion is the same within all subgroups, then observed differences in  $T^2$  must be due to sampling variation alone. In this case, we should pool the information to yield a common estimate, and then apply this estimate to all subgroups. This seems like a plausible expectation in the running example, where the study-to-study variation in effect size is likely to be similar for subgroups *A* and *B*.

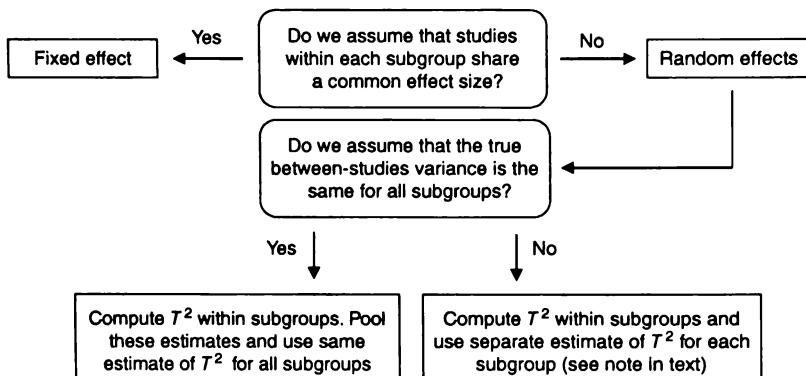
On the other hand, if we anticipate that the true between-studies dispersion may actually differ from one subgroup to the next, then we would estimate  $\tau^2$  within subgroups and use a separate estimate of  $\tau^2$  for each subgroup. For example, suppose that we are assessing an intervention to reduce recidivism among juvenile delinquents, and comparing the effect in subgroups of studies where the delinquents did, or did not, have a history of violence. We might expect to see a wider range of effect sizes in one subgroup than the other.

There is one additional caveat to consider. If we do anticipate that  $\tau^2$  will vary from one subgroup to the next, so that the correct approach is to use separate estimates of  $\tau^2$ , we still need to be sure that there are enough studies within each subgroup to yield an

acceptably accurate estimate of  $\tau^2$ . Generally, if there is only a small number of studies within subgroups, then the estimates of  $\tau^2$  within subgroups are likely to be imprecise. In this case, it makes more sense to use a pooled estimate, since the increased accuracy that we get by pooling more studies is likely to exceed any real differences between groups in the true value of  $\tau^2$ . There is no consensus about what constitutes “a small number of studies” but it would be reasonable to use a number between ten and twenty.

## Summary

The logic outlined above is encapsulated in the flowchart shown in Figure 21.4. If the studies within each subgroup share a common effect size, then we use the fixed-effect model to assign weights to each study (and  $\tau^2$  is zero). Otherwise, we use the random-effects model.



**Figure 21.4** Flowchart for selecting a computational model.

Under random effects we always estimate  $\tau^2$  within subgroups. If we believe that the true value of  $\tau^2$  is the same for all subgroups, then the correct procedure is to pool the estimates obtained within subgroups. If we believe that the true value of  $\tau^2$  varies from one subgroup to the next, the correct procedure is to use a separate estimate for each subgroup. However, if we have only a small number of studies within subgroups (for example, less than twenty) these estimates may be imprecise and therefore it may be preferable to pool the estimates.

## RANDOM EFFECTS WITH SEPARATE ESTIMATES OF $\tau^2$

Here, we proceed through the same set of computations as we did for the fixed-effect model, but this time using random-effects weights, with a separate estimate of  $\tau^2$  for each subgroup.

## Computing the effects

Figure 21.5 is a forest plot of the studies in subgroups A and B. The studies are identical to those in the fixed-effect forest plot (Figure 21.2) but the summary effects, represented by the diamonds, are now based on random-effects weights. The mean effect size for subgroups A and B are 0.325 and 0.610, with variances of 0.006 and 0.004.

Computations are based on the values in Table 21.5. These values are similar to those in Table 21.1, except that the variance for each study now includes the within-study variance and the between-study variance. We did not assume a common value of  $\tau^2$  and therefore used a separate estimate of  $\tau^2$  for each subgroup. In Figure 21.5 this is indicated by the symbols at the right, where we have one value for  $T_A^2$  and another for  $T_B^2$ . In Table 21.5, the column labeled  $T^2$  shows 0.0164 for the A studies and 0.0022 for the B studies.

### *Computations (random effects, separate estimates of $\tau^2$ ) for the A studies*

$$M_A^* = \frac{50.787}{156.512} = 0.3245,$$

$$V_{M_A^*} = \frac{1}{156.512} = 0.0064,$$

$$SE_{M_A^*} = \sqrt{0.0064} = 0.0799,$$

$$LL_{M_A^*} = 0.3245 - 1.96 \times 0.0799 = 0.1678,$$

$$UL_{M_A^*} = 0.3245 + 1.96 \times 0.0799 = 0.4812,$$

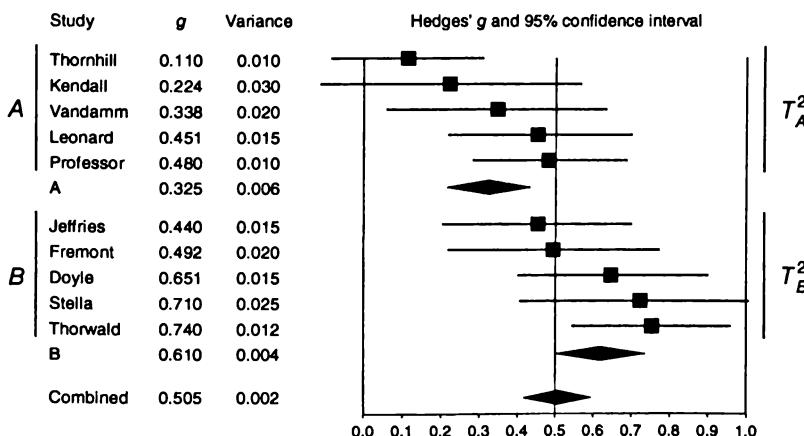


Figure 21.5 Random-effects model (separate estimates of  $\tau^2$ ) – studies and subgroup effects.

**Table 21.5** Random-effects model (separate estimates of  $\tau^2$ ) – computations.

Study	Effect size $Y$	Variance within $V_Y$	Variance between $\tau^2$	Variance total $V$	Weight $W$	Calculated quantities		
						$WY$	$WY^2$	$W^2$
A	Thornhill	0.110	0.0100	0.0164	0.0264	37.846	4.163	1432.308
	Kendall	0.224	0.0300	0.0164	0.0464	21.541	4.825	1.081
	Vandamm	0.338	0.0200	0.0164	0.0364	27.455	9.280	3.137
	Leonard	0.451	0.0150	0.0164	0.0314	31.824	14.353	6.473
	Professor	0.480	0.0100	0.0164	0.0264	37.846	18.166	8.720
	Sum A				156.512	50.787	19.868	5095.179
B	Jefferies	0.440	0.0150	0.0022	0.0172	57.983	25.512	11.225
	Fremont	0.492	0.0200	0.0022	0.0222	44.951	22.116	10.881
	Doyle	0.651	0.0150	0.0022	0.0172	57.983	37.747	24.573
	Stella	0.710	0.0250	0.0022	0.0272	36.702	26.058	18.501
	Thorwald	0.740	0.0120	0.0022	0.0142	70.193	51.943	38.438
	Sum B				267.811	163.376	103.619	15018.633
	Sum				424.323	214.163	123.487	20113.812

$$Z_A^* = \frac{0.3245}{0.0799} = 4.0595,$$

$$p(Z_A^*) < 0.0001,$$

and

$$Q_A^* = 19.868 - \left( \frac{50.787^2}{156.512} \right) = 3.3882. \quad (21.19)$$

Note. The  $Q^*$  statistic computed here, using random-effects weights, is used *only* for the analysis of variance, to partition  $Q^*$  into its various components. Therefore, we do not show a p-value for  $Q^*$ . Rather, the  $Q$  statistic computed using fixed-effect weights (Table 21.2) is the one that reflects the between-studies dispersion, provides a test of homogeneity for the studies within subgroup A, and is used to estimate  $\tau^2$ .

#### Computations (random effects, separate estimates of $\tau^2$ ) for the B studies

$$M_B^* = \frac{163.376}{267.811} = 0.6100,$$

$$V_{M_B^*} = \frac{1}{267.811} = 0.0037,$$

$$SE_{M_B^*} = \sqrt{0.0037} = 0.0611,$$

$$LL_{M_B^*} = 0.6100 - 1.96 \times 0.0611 = 0.4903,$$

$$UL_{M_B^*} = 0.6100 + 1.96 \times 0.0611 = 0.7298,$$

$$Z_B^* = \frac{0.6100}{0.0611} = 9.9833,$$

$$p(Z_B^*) < 0.0001$$

and

$$Q_B^* = 103.619 - \left( \frac{163.376^2}{267.811} \right) = 3.9523. \quad (21.20)$$

### **Computations (random effects, separate estimates of $\tau^2$ ) for all ten studies**

The statistics here are computed using the same value of  $T^2$  as was used within groups (in this case, *not* pooled).

$$M^* = \frac{214.163}{424.323} = 0.5047, \quad (21.21)$$

$$V_{M^*} = \frac{1}{424.323} = 0.0024, \quad (21.22)$$

$$SE_{M^*} = \sqrt{0.0024} = 0.0485,$$

$$LL_{M^*} = 0.5047 - 1.96 \times 0.0485 = 0.4096,$$

$$UL_{M^*} = 0.5047 + 1.96 \times 0.0485 = 0.5999,$$

$$Z^* = \frac{0.5047}{0.0485} = 10.3967,$$

$$p(Z^*) < 0.0001,$$

and

$$Q^* = 123.487 - \left( \frac{214.163^2}{424.323} \right) = 15.3952. \quad (21.23)$$

Statistics (random-effects) are summarized in Table 21.6.

**Table 21.6 Random-effects model (separate estimates of  $\tau^2$ ) – summary statistics.**

	A	B	Combined
$Y$	0.3245	0.6100	0.5047
$V$	0.0064	0.0037	0.0024
$SE_Y$	0.0799	0.0611	0.0485
$LL_Y$	0.1678	0.4903	0.4096
$UL_Y$	0.4812	0.7298	0.5999
$Z$	4.0595	9.9833	10.3967
$p2$	0.0000	0.0000	0.0000
$Q$	3.3882	3.9523	15.3952

### **Comparing the effects**

If we return to Figure 21.5 and excerpt the diamonds for the two subgroups we get Figure 21.6.

The mean effect size for subgroups A and B are 0.325 and 0.610, with variances of 0.006 and 0.004.

Our goal is to compare these two mean effects, and there are several ways that we can proceed. These approaches are algebraically equivalent, and (it follows) yield the

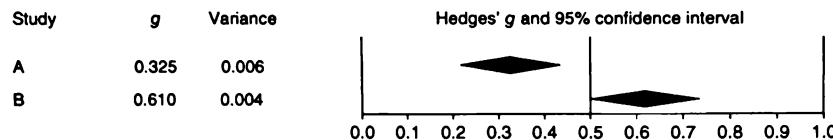


Figure 21.6 Random-effects model (separate estimates of  $\tau^2$ ) – subgroup effects.

same  $p$ -value. Our goal in presenting several approaches is to provide insight into the process.

### Comparing A versus B: a Z-test (Method 1)

We can use a simple Z-test to compare the mean effect for subgroups  $A$  versus  $B$ . The formulas are identical to those used earlier, but we change two symbols to reflect the random-effects model. First, we use a (\*) to indicate that the statistics are based on random-effects weights rather than fixed-effect weights. Second, the null hypothesis is framed as  $\mu_A = \mu_B$ , reflecting the fact that these are mean values, rather than  $\theta_A = \theta_B$ , which we used to refer to common values when we were working with the fixed-effect model.

Let  $\mu_A$  and  $\mu_B$  be the true mean effects underlying subgroups  $A$  and  $B$ , let  $M_A^*$  and  $M_B^*$  be the estimated effects, and let  $V_{M_A^*}$  and  $V_{M_B^*}$  be their variances. If we use  $Diff^*$  to refer to the difference between the two effects and elect to subtract the mean of  $A$  from the mean of  $B$ ,

$$Diff^* = M_B^* - M_A^*. \quad (21.24)$$

The test statistic to compare the two effects is

$$Z_{Diff}^* = \frac{Diff^*}{SE_{Diff^*}}, \quad (21.25)$$

where

$$SE_{Diff^*} = \sqrt{V_{M_A^*} + V_{M_B^*}}. \quad (21.26)$$

Under the null hypothesis that the true mean effect size  $\mu$  is the same for both groups,

$$H_0^* : \mu_A^* = \mu_B^*, \quad (21.27)$$

$Z_{Diff}^*$  would follow the normal distribution. For a two-tailed test the  $p$ -value is given by

$$p^* = 2[1 - (\Phi(|Z_{Diff}^*|))], \quad (21.28)$$

where  $\Phi(Z)$  is the standard normal cumulative distribution.

In the running example

$$Diff^* = 0.6100 - 0.3245 = 0.2856,$$

$$SE_{Diff^*} = \sqrt{0.0064 + 0.0037} = 0.1006,$$

and

$$Z_{Diff} = \frac{0.2856}{0.1006} = 2.8381.$$

The two-tailed  $p$ -value corresponding to  $Z_{Diff}^* = 2.8381$  is 0.0045. This tells us that the mean treatment effect is probably not the same for the  $A$  studies as for the  $B$  studies. In Excel, the function to compute a 2-tailed  $p$ -value for  $Z$  is  $= (1-(NORMSDIST(ABS(Z))))*2$ . Here,  $= (1-(NORMSDIST(ABS(2.8381))))*2$  will return the value 0.0045.

### Comparing $A$ with $B$ : a $Q$ -test based on analysis of variance (Method 2)

We use the same formulas as we did for method 2 under the fixed-effect model, but now apply random-effects weights. Note that this approach only works if we use the same weights to compute the overall effect as we do to compute the effects within groups. In Table 21.5, studies from subgroup  $A$  use the  $T^2$  value of 0.0164 both for computing the subgroup mean and for computing the overall mean. Similarly, studies from subgroup  $B$  use the  $T^2$  value of 0.0022 both for computing the subgroup mean and for computing the overall mean.

We compute the following quantities (where  $SS$  is the sum of squared deviations).

- $Q_A^*$ , the weighted  $SS$  of all  $A$  studies about the mean of  $A$ .
- $Q_B^*$ , the weighted  $SS$  of all  $B$  studies about the mean of  $B$ .
- $Q_{within}^*$ , the sum of  $Q_A^*$  and  $Q_B^*$ .
- $Q_{bet}^*$ , the weighted  $SS$  of the subgroup means about the grand mean.
- $Q^*$ , the weighted  $SS$  of all effects about the grand mean.

We may write  $Q_{within}^* = Q_A^* + Q_B^*$ , to represent the sum of within-group weighted  $SS$ , or more generally, for  $p$  subgroups,

$$Q_{within}^* = \sum_{j=1}^p Q_j^*. \quad (21.29)$$

In the running example

$$Q_{within}^* = 3.3882 + 3.9523 = 7.3406. \quad (21.30)$$

The weighted  $SS$  are additive, such that  $Q^* = Q_{within}^* + Q_{bet}^*$ . Therefore,  $Q_{bet}^*$  can be computed as

$$Q_{bet}^* = Q^* - Q_{within}^*. \quad (21.31)$$

Under the null hypothesis that the effect sizes  $\mu$  are the same for all groups, 1 to  $p$ ,  $Q_{bet}^*$  would be distributed as chi-squared with degrees of freedom equal to  $p - 1$ .

In the running example

$$Q_{bet}^* = 15.3952 - 7.3406 = 8.0547. \quad (21.32)$$

Results are summarized in Table 21.7. Note that the only  $Q$  statistic that we interpret here is the one *between groups*. In the running example, the *Between* line tells us that

**Table 21.7** Random-effects model (separate estimates of  $\tau^2$ ) – ANOVA table.

	$Q^*$	$df$	$p$	Formula
A	3.3882			19.19
B	3.9523			19.20
Within	7.3406			19.30
Between	8.0547	1.0	0.0045	19.32
Total	15.3952			19.23

the difference between groups (the combined effect for A versus B) is statistically significant ( $Q_{bet}^* = 8.0547$ ,  $df = 1$ ,  $p = 0.0045$ ), which means that the effect size is related to the frequency of tutoring. In Excel, the function to compute a  $p$ -value for  $Q$  is =CHIDIST( $Q$ ,  $df$ ). For the test of A versus B, =CHIDIST(8.0547,1) returns 0.0045.

To address the statistical significance of the total variance or the variance within groups, we use the statistics reported using the fixed-effect weights (see Table 21.3) rather than using  $Q^*$  (total),  $Q_A^*$ ,  $Q_B^*$  or  $Q_{within}^*$ .

### Comparing A versus B: a Q-test for heterogeneity (Method 3)

Finally, we could treat the subgroups as if they were studies and perform a test for heterogeneity across studies. If we extract the two subgroup lines and the total line from Figure 21.5 and replace the diamonds with squares we get Figure 21.7.

Concretely, we start with two *studies* with effect sizes of 0.324 and 0.610, and variances of 0.006 and 0.004. Then, we apply the usual meta-analysis methods to compute  $Q$ . Concretely, using the values in Table 21.8, and applying (11.2) and subsequent formulas, we compute

$$M^* = \frac{214.163}{424.323} = 0.5047 \quad (21.33)$$

and

$$V_M^* = \frac{1}{424.323} = 0.0024, \quad (21.34)$$

$$Q = 116.146 - \left( \frac{214.163^2}{424.323} \right) = 8.0547,$$

$$df = 2 - 1 = 1,$$

**Figure 21.7** Random-effects model (separate estimates of  $\tau^2$ ) – treating subgroups as studies.

**Table 21.8** Random-effects model (separate estimates of  $\tau^2$ ) – subgroups as studies.

Study	Effect size $Y$	Variance within $V_Y$	Variance between $T^2$	Variance total $V$	Weight $W$	Calculated quantities		
						$WY$	$WY^2$	$W^2$
A	0.3245	0.0064	0.0000	0.0064	156.512	50.787	16.480	24495.944
B	0.6100	0.0037	0.0000	0.0037	267.811	163.376	99.666	71722.774
					424.323	214.163	116.146	96218.718

and

$$p(Q = 8.0547, df = 1) = 0.0045,$$

where  $Q$  represents the weighted sum of squares for Studies A and B about the grand mean. For  $Q = 8.0547$  and  $df = 1$ , the  $p$ -value is 0.0045.

In Excel, the function to compute a  $p$ -value for  $Q$  is =CHIDIST( $Q, df$ ). For the test of A versus B, =CHIDIST(8.0547,1) returns 0.0045.

#### Quantify the magnitude of the difference

The difference and confidence interval are given by (21.17) and (21.18):

$$Diff^* = 0.6100 - 0.3245 = 0.2856,$$

$$SE_{Diff^*} = \sqrt{0.0064 + 0.0037} = 0.1006,$$

$$LL_{Diff^*} = 0.2856 - 1.96 \times 0.1006 = 0.0883,$$

and

$$UL_{Diff^*} = 0.2856 + 1.96 \times 0.1006 = 0.4828.$$

In words, the true difference between the effect in the A studies, as opposed to the B studies, probably falls in the range of 0.09 to 0.48.

#### RANDOM EFFECTS WITH POOLED ESTIMATE OF $\tau^2$

Here, we show the computation of summary effects within subgroups, using a random-effects model with a pooled estimate of  $\tau^2$ , which we refer to as  $\tau^2_{within}$ . We illustrate the procedure in Figure 21.8. Note the common value of  $\tau^2_{within}$  is assumed to apply to both subgroups.

#### Formula for estimating a pooled $\tau^2$

To estimate the pooled  $\tau^2$ , proceed as follows. Recall (12.2) to (12.5) that to estimate  $\tau^2$  for a single collection of studies we use

$$T^2 = \frac{Q - df}{C}, \quad (21.35)$$

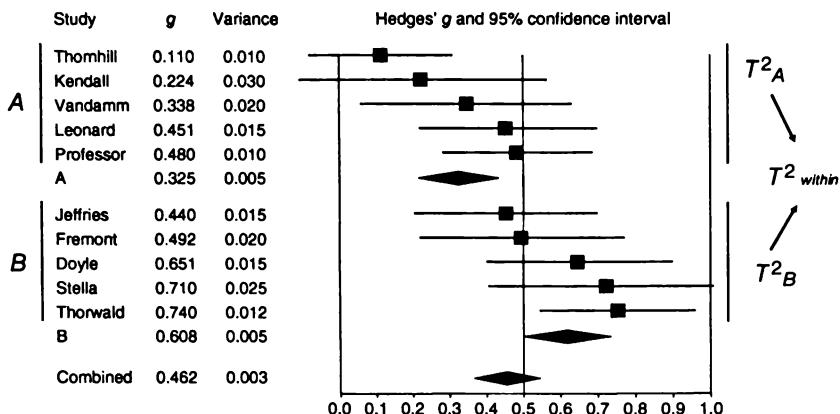


Figure 21.8 Random-effects model (pooled estimate of  $\tau^2$ ) – studies and subgroup effects.

where

$$df = k - 1, \quad (21.36)$$

where  $k$  is the number of studies, and

$$C = \sum W_i - \frac{\sum W_i^2}{\sum W_i}, \quad (21.37)$$

In these equations,  $Q - df$  is the excess (observed minus expected) sum of squared deviations from the weighted mean, and  $C$  is a scaling factor.

Similarly, to yield a pooled estimate of  $\tau^2$  we sum each element ( $Q$ ,  $df$ , and  $C$ ) across subgroups and then perform the same computation. Concretely,

$$T_{within}^2 = \frac{\sum_{j=1}^p Q_j - \sum_{j=1}^p df_j}{\sum_{j=1}^p C_j}. \quad (21.38)$$

While the true value of  $\tau_{within}^*$  cannot be less than zero (a variance cannot be negative), this method of estimating  $\tau_{within}^*$  can yield a negative value due to sampling issues (when the observed dispersion is less than we would expect by chance). In this case, the estimate  $T_{within}^*$  is set to zero.

### Computing the effects

Subgroup A yielded an estimate of 0.0164 while subgroup B yielded an estimate of 0.0122, represented in Figure 21.6 as  $T_A^2$  and  $T_B^2$ . We will pool these two estimates to yield a pooled value, represented as  $T_{within}^2$ , of 0.0097 (see (21.39)). This is the value used to assign weights in Table 21.10.

**Table 21.9** Statistics for computing a pooled estimate of  $\tau^2$ .

Group	<i>Q</i>	<i>df</i>	<i>C</i>
A	8.4316	4	269.8413
B	4.5429	4	241.6667
Sum	12.9745	8	511.5079

In the running example, the values within each group were computed earlier for *A* and *B*. Table 21.9 shows the values needed to calculate a pooled estimate  $T_{within}^2$  for the running example.

Then,

$$T_{within}^2 = \frac{12.9745 - 8}{511.508} = 0.00974. \quad (21.39)$$

Computations below are based on the values in Table 21.10. These are similar to Table 21.5, except that we now assume that all groups have the same  $\tau^2$ , and use a common estimate. In Table 21.10 the same estimate of  $\tau^2$  (0.0097) is applied to all ten studies.

#### *Computations (random effects, pooled estimate of $\tau^2$ ) for the A studies*

$$M_A^* = \frac{65.161}{200.652} = 0.3247,$$

$$V_{M_A^*} = \frac{1}{200.652} = 0.0050,$$

$$SE_{M_A^*} = \sqrt{0.0050} = 0.0706,$$

**Table 21.10** Random-effects model (pooled estimate of  $\tau^2$ ) – computations.

Study	Effect size $\gamma$	Variance within $V_\gamma$	Variance between $\tau^2$	Variance total $V$	Weight $W$	Calculated quantities			
						$WY$	$WY^2$	$W^2$	
A	Thornhill	0.110	0.0100	0.0097	0.0197	50.697	5.577	0.613	2570.150
	Kendall	0.224	0.0300	0.0097	0.0397	25.173	5.639	1.263	633.678
	Vandamm	0.338	0.0200	0.0097	0.0297	33.642	11.371	3.843	1131.752
	Leonard	0.451	0.0150	0.0097	0.0247	40.445	18.241	8.226	1635.767
	Professor	0.480	0.0100	0.0097	0.0197	50.697	24.334	11.681	2570.150
	Sum A				200.652	65.161	25.627	8541.498	
B	Jefferies	0.440	0.0150	0.0097	0.0247	40.445	17.796	7.830	1635.767
	Fremont	0.492	0.0200	0.0097	0.0297	33.642	16.552	8.143	1131.752
	Doyle	0.651	0.0150	0.0097	0.0247	40.445	26.329	17.140	1635.767
	Stella	0.710	0.0250	0.0097	0.0347	28.798	20.446	14.517	829.299
	Thorwald	0.740	0.0120	0.0097	0.0217	46.030	34.062	25.206	2118.721
	Sum B				189.358	115.185	72.837	7351.306	
	Sum				390.010	180.346	98.463	15892.804	

$$LL_{M_A^*} = 0.3247 - 1.96 \times 0.0706 = 0.1864,$$

$$UL_{M_A^*} = 0.3247 + 1.96 \times 0.0706 = 0.4631,$$

$$Z_A^* = \frac{0.3247}{0.0706} = 4.6601,$$

$$p(Z_A^*) < 0.0001,$$

and

$$Q_A^* = 25.627 - \left( \frac{65.161^2}{200.652} \right) = 4.4660. \quad (21.40)$$

Note. The  $Q^*$  statistic computed here, using random-effects weights, is used *only* for the analysis of variance, to partition  $Q^*$  into its various components. Therefore, we do not show a  $p$ -value for  $Q^*$ . Rather, the  $Q$  statistic computed using fixed-effect weights (above) is the one that reflects the between-studies dispersion, provides a test of homogeneity for the studies within subgroup  $A$ , and is used to estimate  $\tau_{within}^2$ .

#### **Computations (random effects, pooled estimate of $\tau^2$ ) for the B studies**

$$M_B^* = \frac{115.185}{189.358} = 0.6083,$$

$$V_{M_B^*} = \frac{1}{189.358} = 0.0053,$$

$$SE_{M_B^*} = \sqrt{0.0053} = 0.0727,$$

$$LL_{M_B^*} = 0.6083 - 1.96 \times 0.0727 = 0.4659,$$

$$UL_{M_B^*} = 0.6083 + 1.96 \times 0.0727 = 0.7507,$$

$$Z_B^* = \frac{0.6083}{0.0727} = 8.3705,$$

$$p(Z_B^*) < 0.0001,$$

and

$$Q_B^* = 72.837 - \left( \frac{115.185^2}{189.358} \right) = 2.7706. \quad (21.41)$$

#### **Computations (random effects, pooled estimate of $\tau^2$ ) for all ten studies**

The statistics here are computed using the same value of  $T^2$  as was used within groups (in this case, the pooled estimate,  $T_{within}^2$ ).

$$M^* = \frac{180.346}{390.010} = 0.4624, \quad (21.42)$$

$$V_{M^*} = \frac{1}{390.010} = 0.0026, \quad (21.43)$$

$$SE_{M^*} = \sqrt{0.0026} = 0.0506,$$

$$LL_{M^*} = 0.4624 - 1.96 \times 0.0506 = 0.3632,$$

$$UL_{M^*} = 0.4624 + 1.96 \times 0.0506 = 0.5617,$$

$$Z^* = \frac{0.4624}{0.0506} = 9.1321,$$

$$p(Z^*) < 0.0001,$$

and

$$Q^* = 98.463 - \left( \frac{180.346^2}{390.010} \right) = 15.0690. \quad (21.44)$$

The statistics computed above are summarized in Table 21.11.

### Comparing the effects

If we return to Figure 21.8 and excerpt the diamonds for the two subgroups we get Figure 21.9. The mean effect size for subgroups *A* and *B* are 0.325 and 0.608, with variances of 0.005 and 0.005.

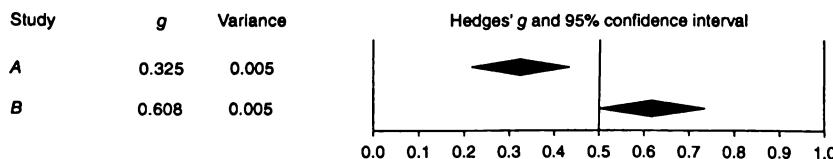
Our goal is to compare these two mean effects, and there are several ways that we can proceed. These approaches are algebraically equivalent, and (it follows) yield the same *p*-value.

### Comparing *A* versus *B*: a Z-test (Method 1)

We can use a Z-test to compare the mean effect for subgroups *A* versus *B*. The null hypothesis and formulas are the same as those for the prior case (where we did not

**Table 21.11** Random-effects model (pooled estimate of  $\tau^2$ ) – summary statistics.

	A	B	Combined
<i>Y</i>	0.3247	0.6083	0.4624
<i>V</i>	0.0050	0.0053	0.0026
<i>SE<sub>Y</sub></i>	0.0706	0.0727	0.0506
<i>LL<sub>Y</sub></i>	0.1864	0.4659	0.3632
<i>UL<sub>Y</sub></i>	0.4631	0.7507	0.5617
<i>Z</i>	4.6001	8.3705	9.1321
<i>p</i> <sup>2</sup>	0.0000	0.0000	0.0000
<i>Q</i>	4.4660	2.7706	15.0690



**Figure 21.9** Random-effects model (pooled estimate of  $\tau^2$ ) – subgroup effects.

assume a common value for  $\tau^2$ ). If we elect to subtract the mean of A from the mean of B,

$$Diff^* = M_B^* - M_A^*, \quad (21.45)$$

the test statistic to compare the two effects is

$$Z_{Diff}^* = \frac{Diff^*}{SE_{Diff}}, \quad (21.46)$$

where

$$SE_{Diff} = \sqrt{V_{M_A^*} + V_{M_B^*}}. \quad (21.47)$$

Under the null hypothesis that the true mean effect size  $\mu_i$  is the same for both groups,

$$H_0 : \mu_A = \mu_B. \quad (21.48)$$

$Z_{Diff}^*$  would follow the normal distribution. For a two-tailed test the  $p$ -value is given by

$$p^* = 2[1 - (\Phi(|Z_{Diff}^*|))], \quad (21.49)$$

where  $\Phi(Z)$  is the standard normal cumulative distribution.

In the running example

$$Diff^* = 0.6083 - 0.3247 = 0.2835,$$

$$SE_{Diff} = \sqrt{0.0050 + 0.0053} = 0.1013,$$

and

$$Z_{Diff}^* = \frac{0.2835}{0.1013} = 2.7986.$$

The two-tailed  $p$ -value corresponding to  $Z_{Diff}^* = 2.7986$  is 0.0051. This tells us that the mean effect is probably not the same for the A studies as for the B studies. In Excel, the function to compute a 2-tailed  $p$ -value for Z is  $=(1-(NORMSDIST(ABS(Z))))^*2$ . Here,  $=(1-(NORMSDIST(ABS(2.7986))))^*2$  will return the value 0.0045.

### Comparing A with B: a Q-test based on analysis of variance (Method 2)

Again, we apply the same formulas as we did for the prior case, but this time using the random-effects weights based on a pooled estimate of  $\tau^2$ . Note that this approach only works if we use the same weights to compute the overall effect as we do to compute the effects within groups. In Table 21.10 we used a  $\tau^2$  value of 0.0097 for all ten studies, and this is the value used to sum *within* subgroups and also to sum across subgroups.

We compute the following quantities (where SS is the sum of squared deviations).

- $Q_A^*$ , the weighted SS of all A studies about the mean of A.
- $Q_B^*$ , the weighted SS of all B studies about the mean of B.
- $Q_{within}^*$ , the sum of  $Q_A^*$  and  $Q_B^*$ .
- $Q_{bet}^*$ , the weighted SS of the subgroup means about the grand mean.
- $Q^*$ , the weighted SS of all effects about the grand mean.

**Table 21.12** Random-effects model (pooled estimate of  $\tau^2$ ) – ANOVA table.

	$Q^*$	$df$	$p$	Formula
A	4.4660			19.40
B	2.7706			19.41
Within	7.2366			19.51
Between	7.8324	1	0.0051	19.53
Total	15.0690			19.44

We may write  $Q_{within}^* = Q_A^* + Q_B^*$ , to represent the sum of within-group weighted SS, or more generally, for  $p$  subgroups,

$$Q_{within}^* = \sum_{j=1}^p Q_j^*. \quad (21.50)$$

In the running example (Table 21.12),

$$Q_{within}^* = 4.4660 + 2.7706 = 7.2366. \quad (21.51)$$

The weighted SS are additive, such that  $Q^* = Q_{within}^* + Q_{bet}^*$ . Therefore,  $Q_{bet}^*$  can be computed as

$$Q_{bet}^* = Q^* - Q_{within}^*. \quad (21.52)$$

Under the null hypothesis that the true mean effect size  $\mu$  is the same for all groups, 1 to  $p$ ,  $Q_{bet}^*$  would be distributed as chi-squared with degrees of freedom equal to  $p - 1$ .

In the running example

$$Q_{bet}^* = 15.0690 - 7.2366 = 7.8324. \quad (21.53)$$

The only  $Q$  statistic that we interpret here is the one between groups. In the running example, the *Between* line tells us that the difference between groups (the combined effect for A versus B) is statistically significant ( $Q_{bet}^* = 7.8324$   $df = 1$ ,  $p = 0.0051$ ), which means that the effect size is related to the frequency of tutoring. In Excel, the function to compute a  $p$ -value for  $Q$  is =CHIDIST( $Q$ ,  $df$ ). For the test of A versus B, =CHIDIST(7.8324, 1) returns 0.0051.

To address the statistical significance of the total variance or the variance within groups, we use the statistics reported using the fixed-effect weights (Table 21.3) rather than using  $Q_{total}^*$ ,  $Q_A^*$ ,  $Q_B^*$  or  $Q_{within}^*$ .

### Comparing A versus B: a Q-test for heterogeneity (Method 3)

Finally, we could treat the subgroups as if they were studies and perform a test for heterogeneity across studies. If we extract the two subgroup lines and the total line from Figure 21.8 and replace the diamonds with squares we obtain Figure 21.10.

Concretely, we start with two *studies* with effect sizes of 0.325 and 0.608, and variances of 0.005 and 0.005. Then, we apply the usual meta-analysis methods to

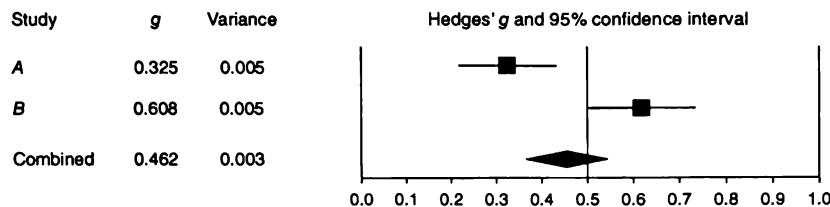


Figure 21.10 Random-effects model (pooled estimate of  $\tau^2$ ) – treating subgroups as studies.

Table 21.13 Random-effects model (pooled estimate of  $\tau^2$ ) – subgroups as studies.

Study	Effect size $Y$	Variance within $V_Y$	Variance between $\tau^2$	Variance total $V$	Weight $W$	Calculated quantities		
						$WY$	$WY^2$	$W^2$
A	0.3247	0.0050	0.0000	0.0050	200.652	65.161	21.161	40261.386
B	0.6083	0.0053	0.0000	0.0053	189.358	115.185	70.066	35856.405

compute  $Q$ . Concretely, using the values in Table 21.13, and applying (11.2) and subsequent formulas, we compute

$$M^* = \frac{180.346}{390.010} = 0.4624, \quad (21.54)$$

$$V_M^* = \frac{1}{390.010} = 0.0026, \quad (21.55)$$

$$Q = 91.227 - \left( \frac{180.346^2}{390.010} \right) = 7.8324,$$

$$df = 2 - 1 = 1,$$

and

$$p(Q = 7.8324, df = 1) = 0.0051,$$

where  $Q$  represents the weighted sum of squares for Studies A and B about the grand mean. For  $Q = 7.8324$  and  $df = 1$ , the  $p$ -value is 0.0051.

In Excel, the function to compute a  $p$ -value for  $Q$  is  $=\text{CHIDIST}(Q, df)$ . For example,  $=\text{CHIDIST}(7.8324, 1)$  returns 0.0051.

### Quantify the magnitude of the difference

The difference and confidence interval are given by (21.17) and (21.18):

$$Diff^* = 0.6083 - 0.3247 = 0.2835,$$

$$SE_{Diff^*} = \sqrt{0.0050 + 0.0053} = 0.1013,$$

$$LL_{Diff^*} = 0.2835 - 1.96 \times 0.1013 = 0.0850,$$

and

$$UL_{Diff} = 0.2835 + 1.96 \times 0.1013 = 0.4821.$$

In words, the true difference between the effect in the *A* studies, as opposed to the *B* studies, probably falls in the range of 0.09 to 0.48.

## THE PROPORTION OF VARIANCE EXPLAINED

In primary studies, a common approach to describing the impact of a covariate is to report the proportion of variance explained by that covariate. That index,  $R^2$ , is defined as the ratio of explained variance to total variance,

$$R^2 = \frac{\sigma^2_{explained}}{\sigma^2_{total}} \quad (21.56)$$

or, equivalently,

$$R^2 = 1 - \left( \frac{\sigma^2_{unexplained}}{\sigma^2_{total}} \right). \quad (21.57)$$

This index is intuitive because it can be interpreted as a ratio, with a range of 0 to 1 (or expressed as a percentage in the range of 0% to 100%). Many researchers are familiar with this index, and have a sense of what proportion of variance is likely to be explained by different kinds of covariates or interventions.

This index cannot be applied directly to meta-analysis for the following reason. In a primary study, a covariate that explains all of the variation in the dependent variable will reduce the error to zero (and  $R^2$ , the proportion of variance explained, would reach 100%).

For example, Figure 21.11 depicts a primary study with 10 participants. All those in group *A* have the same score (0.3) and all those in group *B* have the same score (0.6). Since the variance *within* each subgroup is 0.0, group membership explains 100% of the original variance, and  $R^2$  is 100%. In a real study, of course, there would be some variance within groups and  $R^2$  would be less than 100%, but the fact that  $R^2$  can potentially reach 100% is part of what makes this index intuitive.

By contrast, consider what happens in a meta-analysis if we have two subgroups of studies. We assume that there are five studies in each subgroup, with a *true* summary effect (say, a standardized mean difference) of 0.30 for each study in subgroup *A* and of 0.70 for each study in subgroup *B*. However, while the true effect is identical for each study within its subgroup, the observed effects will differ from each other because of random error.

Thus, the variance within groups, while smaller than the variance between groups, can never approach zero. If the within-study error is a substantial portion of the total variance observed (say, 75%), then the upper limit of  $R^2$  would be only 25%. As such, two important qualities of the index (the fact that it has a natural scale of 0% to 100% and the fact that it has the same range across studies) would no longer apply.

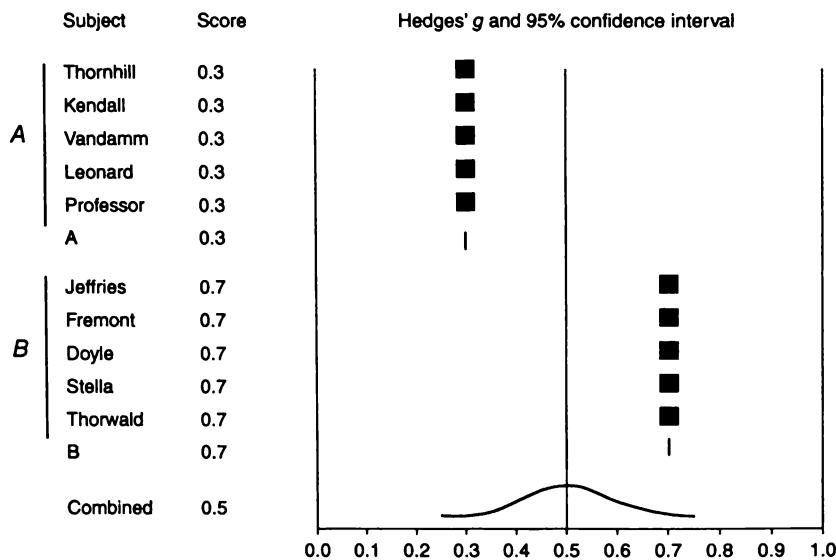


Figure 21.11 A primary study showing subjects within groups.

Since the problem with using  $R^2$  is the fact that study-level covariates in a meta-analysis can address only the true variance  $\tau^2$  (and not the within-study variance  $v$ ), the logical solution is to redefine  $R^2$  (or to define a new index) that is based solely on the true variance. Rather than defining  $R^2$  as the proportion of *total* variance explained by the covariates, we will define it as the proportion of *true* variance explained by the covariates. Since the true variance is estimated as  $T^2$ , this gives us

$$R^2 = \frac{T^2_{\text{explained}}}{T^2_{\text{total}}}, \quad (21.58)$$

or

$$R^2 = 1 - \left( \frac{T^2_{\text{unexplained}}}{T^2_{\text{total}}} \right). \quad (21.59)$$

In the context of subgroups, the numerator in (21.59) is the between-studies variance within subgroups, and the denominator is the total between-studies variance (within-subgroups plus between-subgroups). Therefore, the equation can be written

$$R^2 = 1 - \left( \frac{T^2_{\text{within}}}{T^2_{\text{total}}} \right). \quad (21.60)$$

In the running example,  $T^2$  for the full set of studies was 0.0299 (see page 166), and  $T^2$  computed by working within subgroups and then pooling across subgroups was

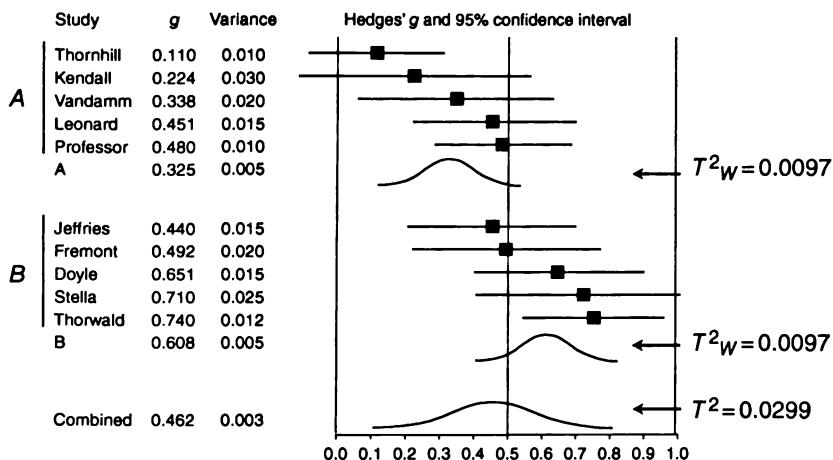


Figure 21.12 Random-effects model – variance within and between subgroups.

0.0097 (see page 183). This gives us

$$R^2 = 1 - \left( \frac{0.0097}{0.0299} \right) = 0.6745. \quad (21.61)$$

In Figure 21.12 we have superimposed a normal curve for the distribution of true effects within each subgroup of studies, and also across all ten studies. The relatively narrow dispersion within groups is based on the  $T^2$  of 0.0097, while the relatively wide dispersion across groups is based on the  $T^2$  of 0.0299, and  $R^2$  captures this change.

The same idea is shown from another perspective in Figure 21.13. On the top line, 34% of the total variance was within studies and 66% was between studies (which is also the definition of  $I^2$ ). The within-studies variance cannot be explained by study level covariates, and so is removed from the equation and we focus on the shaded part. On the bottom line, the type of intervention is able to explain 67% of the *relevant* variance, leaving 33% unexplained. Critically, the 67% and 33% sum to 100%, since we are concerned only with the variance between studies.

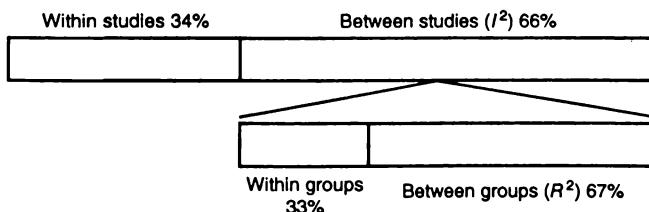


Figure 21.13 Proportion of variance explained by subgroup membership.

Note 1. While the  $R^2$  index has a range of 0 to 1 (0% to 100%) in the population, it is possible for sampling error to yield an observed value of  $R^2$  that falls outside of this range. In that case, the value is set to either 0 (0%) or 1 (100%).

Note 2. The  $R^2$  index only makes sense if we are using a random-effects model, which allows us to think about explaining some of the between-studies variance. Under the fixed-effect model the between-studies variance is set at zero and cannot be changed. Also, the computational model proposed here for estimating  $R^2$  only works for the case where we assume that  $\tau^2$  is the same for all subgroups.

## MIXED-EFFECTS MODEL

In this volume we have been using the term fixed effect to mean that the effect is identical (*fixed*) across all relevant studies (within the full population, or within a subgroup).

In fact the use of the term *fixed effect* in connection with meta-analysis is at odds with the usual meaning of *fixed effects* in statistics. A more suitable term for the fixed-effect meta-analysis might be a *common-effect* meta-analysis. The term *fixed effects* is traditionally used in another context with a different meaning. Concretely, we can talk about the subgroups as being *fixed* in the sense of fixed rather than random. For example, if we want to compare the treatment effect for a subgroup of studies that enrolled only males versus a subgroup of studies that enrolled only females, then we would assume that the subgroups are *fixed* in the sense that anyone who wanted to perform this analysis would need to use these same two subgroups (male and female). By contrast, if we have subgrouping of studies by country, then we might prefer to treat the subgroups as random. A random-effects assumption across subgroups of studies in the US, UK, Japan, Australia and Sweden would allow us to infer what the effect might be in a study in Israel, by assuming it comes from the same random-effects distribution. In this chapter we assume that when we are interested in comparing subgroups we make an assumption of the first type, which means that anyone who performs this comparison must use the same set of subgroups.

We mention this here for two reasons. One is to alert the reader that in the event that the subgroups have been sampled at random from a larger pool (as in the example of countries), then we are able to take this additional source of variability into account. The mechanism for doing so is beyond the scope of an introductory book.

The other reason is to explain the meaning of the term *mixed model*, which is sometimes used to describe subgroup analyses. As explained in this chapter the summary effect *within subgroups* can be computed using either a fixed-effect model or a random-effects model. As outlined immediately above, the difference *across subgroups* can be assessed using either a fixed-effects model or a random-effects model (although the meaning of *fixed* is different here). This leads to the following nomenclature.

If we use a fixed-effect model within subgroups and also across subgroups, the analysis is called a fixed-effects analysis. If we use a random-effects model within subgroups and a fixed-effect model across subgroups (the approach that we generally advocate), the model is called a mixed-effects model. We have the further possibility

of assuming random effects both within and across subgroups; such a model is called a random-effects (or fully random-effects) model.

## OBTAINING AN OVERALL EFFECT IN THE PRESENCE OF SUBGROUPS

In the tables and forest plots presented in this chapter we presented a summary effect for each subgroup and also for the total population. Since our primary concern has been with looking at difference between subgroups we paid little attention to the value for the total population. Here, we consider if that value should be reported at all, and if so, how it should be computed.

### Should we report a summary effect across all subgroups?

The question of whether or not we should report a summary effect across all subgroups depends on our goals and also on the nature of the data.

Suppose the primary goal of the analysis is to see if a treatment is more effective among acute patients than among chronic patients, and it emerges that the treatment is very effective in one group but harmful in the other. In this case, the take-home message should be that we need to look at each group separately. To report that the treatment is moderately effective (on average) would be a bad idea since this is true for neither group and misrepresents the core finding. In this case, it would be better to report the effect for the separate subgroups only.

By contrast, if it turns out that the treatment is equally effective (or nearly so) in both subgroups, then it might be helpful to report a combined effect to serve as a summary. This would probably be the case also if there are minor differences among groups, but the substantive implication of the treatment (or the relationship) is the same for all groups. This is especially true if there are many subgroups, and the reader will be looking for a single number that is easy to recall.

If we do decide to report a combined effect across subgroups, we need to be clear about what this value represents, since this determines how it will be computed. The basic options are explained below.

### Option 1. Combine subgroup means and ignore between-subgroup variance

One option is to compute the weighted mean of the subgroup means. In other words, we treat each subgroup as a study and perform a fixed-effect analysis using the mean effect and variance for each subgroup. In this chapter, we showed three versions of this approach.

These computations were shown for the fixed-effect model in (21.14) and (21.15) and where we computed the weighted mean of the two subgroups. Note that we would get the identical values if we worked with the original ten studies and weighted each by its fixed-effect weight (see (21.3) and (21.4)).

These computations were shown for the random-effects model with separate estimates of  $\tau^2$  in (21.33) and (21.34), where we computed the weighted mean of the two

subgroups. Note that we would get the identical values if we worked with the original ten studies and weighted each by its random-effects weight, with a separate estimate of  $\tau^2$  for each subgroup (see (21.21) and (21.22)).

These computations were shown for the random-effects model with a pooled estimate of  $\tau^2$  in (21.54) and (21.55), where we computed the weighted mean of the two subgroups. Note that we would get the identical values if we worked with the original ten studies and weighted each by its random-effects weight, with a pooled estimate of  $\tau^2$  (see (21.42) and (21.43)).

In all three cases, the combined effect refers to no actual population but is rather the average of two different populations. If the subgroups were male and female then the combined effect is the expected effect in a population that included both males and females (in the same proportions as in the subgroups). As always, the standard error of the mean speaks to the precision of the mean, and not to the dispersion of effects across subgroups (which is treated as zero).

### **Option 2. Combine subgroup means, and model the between-subgroup variance**

A second option is to assume a random-effects model across subgroups. In other words, all the formulas and concepts discussed in Chapter 12 are applied here, except that the unit of analysis is the subgroup rather than the study. This would make sense if the subgroups have been sampled at random from a larger group of relevant subgroups. For example, we have the mean effect of a treatment in the US and in Australia, but we want to estimate what the mean effect of that treatment would be across all relevant countries.

In this case we need to address precisely the same kinds of issues we addressed when discussing heterogeneity in Chapter 12. First, we compute a measure of between-subgroups dispersion,  $T^2_{het}$ . Then, we compute a weighted mean of the subgroups, where the weights are based on the within-subgroup error and the between-subgroups variance. To the extent that the subgroup means differ from each other, the standard error of the combined effect will be increased (but this additional error will be diminished as additional subgroups are added).

We can also focus on the dispersion itself (as in Chapters 16 and 17). For example, we can use the estimate of  $\tau^2_{het}$  to build a prediction interval that gives us the expected range of effect sizes for the next country (in our example) selected at random.

### **Option 3. Perform a separate random-effects analysis on the full set of studies.**

If we want to report a combined effect across subgroups, then a third option is simply to perform a separate random-effects meta-analysis including all of the studies, and ignoring subgroup membership. Rather than estimate  $\tau^2$  within subgroups (as we did before) we estimate it across all studies, and so it will tend to be larger.

## Comparing the options

When our primary goal is to assess differences among subgroups, and use an analysis of variance table as part of the process, the combined effects across subgroups are computed using option 1. This yields a set of internally consistent data.

If we really care about the combined effect across subgroups then options 2 and 3 are the more logical choices. If the subgroups really have been selected at random from a larger set, then option 2 allows us to model the different sources of error separately and obtain a better estimate of the true confidence interval for the combined effect (as well as discuss prediction intervals for a future subgroup), and is probably the better choice. This assumes, of course, that we have sufficient information to obtain a reasonably precise estimate of the variance among subgroups. By contrast, if the subgrouping is not of major importance, or if multiple different subgroupings of the studies are being considered, then option 3 is the more logical choice.

### SUMMARY POINTS

- Just as we can use *t*-tests or analysis of variance in primary studies to assess the relationship between group membership and outcome, we can use analogs of these procedures in meta-analysis to assess the relationship between subgroup membership and effect size.
- We presented three methods that can be used to compare the mean effect across subgroups. To compare the mean effect in two groups we can use a Z-test. To compare the mean effect in two or more groups we can use analysis of variance (modified for use with subgroups) or the *Q*-test of homogeneity. All three procedures are mathematically equivalent.
- These analyses may be performed using either the fixed-effect or the random-effects model within groups, but in most cases the latter is appropriate.
- In primary studies we use  $R^2$  to reflect the proportion of variance explained by group membership. An analogous index, which reflects the proportion of true variance explained by subgroup membership, can be used for meta-analysis.



# Meta-Regression

---

### Introduction

Fixed-effect model

Fixed or random effects for unexplained heterogeneity

Random-effects model

---

## INTRODUCTION

In primary studies we use regression, or multiple regression, to assess the relationship between one or more covariates (moderators) and a dependent variable. Essentially the same approach can be used with meta-analysis, except that the covariates are at the level of the study rather than the level of the subject, and the dependent variable is the effect size in the studies rather than subject scores. We use the term *meta-regression* to refer to these procedures when they are used in a meta-analysis.

The differences that we need to address as we move from primary studies to meta-analysis for regression are similar to those we needed to address as we moved from primary studies to meta-analysis for subgroup analyses. These include the need to assign a weight to each study and the need to select the appropriate model (fixed versus random effects). Also, as was true for subgroup analyses, the  $R^2$  index, which is used to quantify the proportion of variance explained by the covariates, must be modified for use in meta-analysis.

With these modifications, however, the full arsenal of procedures that fall under the heading of multiple regression becomes available to the meta-analyst. We can work with sets of covariates, such as three variables that together define a treatment, or that allow for a nonlinear relationship between covariates and the effect size. We can enter covariates into the analysis using a pre-defined sequence and assess the impact of any set, over and above the impact of prior sets, to control for confounding variables. We can incorporate both categorical (for example, dummy-coded) and continuous variables as covariates. We can use these procedures both to assess the impact of covariates and also to predict the effect size in studies with specific characteristics.

Multiple regression incorporates a wide array of procedures, and we cannot cover these fully in this volume. Rather, we assume that the reader is familiar with multiple regression in primary studies, and our goal here is to show how the same techniques used in primary studies can be applied to meta-regression.

As is true in primary studies, where we need an appropriately large ratio of *subjects* to covariates in order for the analysis be to meaningful, in meta-analysis we need an appropriately large ratio of *studies* to covariates. Therefore, the use of meta-regression, especially with multiple covariates, is not a recommended option when the number of studies is small. In primary studies some have recommended a ratio of at least ten subjects for each covariate, which would correspond to ten studies for each covariate in meta-regression. In fact, though, there are no hard and fast rules in either case.

### FIXED-EFFECT MODEL

As we did when discussing subgroup analysis, we start with the fixed-effect model, which is simpler, and then move on to the random-effects model, which is generally more appropriate.

#### The BCG data set

Various researchers have published studies that assessed the impact of a vaccine, known as BCG, to prevent the development of tuberculosis (TB). With the re-emergence of TB in the United States in recent years (including many drug resistant cases), researchers needed to determine whether or not the BCG vaccine should be recommended. For that reason, Colditz *et al.* (1994) reported a meta-analysis of these studies, and Berkey *et al.* (1995) showed how meta-regression could be used in an attempt to explain some of the variance in treatment effects.

The forest plot is shown in Figure 22.1. The effect size is the risk ratio, with a risk ratio of 0.10 indicating that the vaccine reduced the risk of TB by 90%, a risk ratio of 1.0 indicating no effect, and risk ratios higher than 1.0 indicating that the vaccine increased the risk of TB. Studies are sorted from most effective to least effective. As always for an analysis of risk ratios, the analysis was performed using log transformed values and then converted back to risk ratios for presentation.

Using a fixed-effect analysis the risk ratio for the 13 studies is 0.650 with a confidence interval of 0.601 to 0.704, which says that the vaccine decreased the risk of TB by at least 30% and possibly by as much as 40%. In log units the risk ratio is -0.430 with a standard error of 0.040. The Z-value is -10.625 ( $p < 0.0001$ ), which allows us to reject the null hypothesis of no effect.

At least as interesting, however, is the substantial variation in the treatment effects, which ranged from a risk ratio of 0.20 (an 80% reduction in risk) to 1.56 (a 56% *increase* in risk). While some of this variation is due to within-study error, some of it reflects variation in the true effects. The  $Q$ -statistic is 152.233 with  $df = 12$  and  $p < 0.0001$ .  $T^2$  is 0.309, which yields a prediction interval of approximately 0.14 to

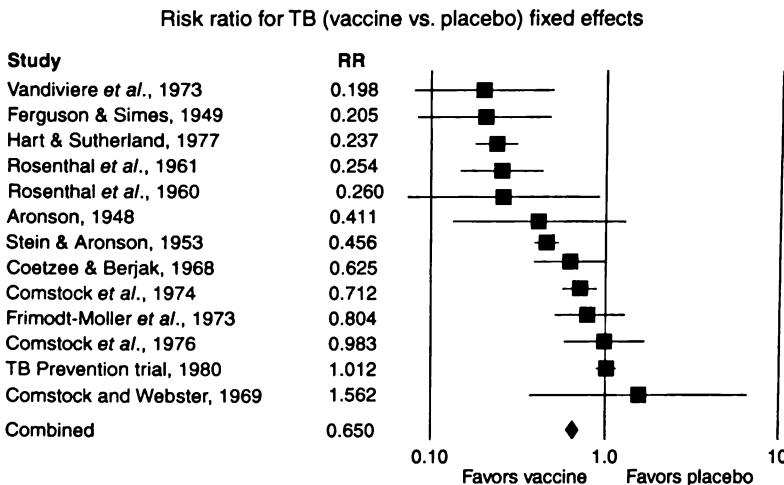


Figure 22.1 Fixed-effect model – forest plot for the BCG data.

1.77. If the effects are normally distributed in log units, the true risk ratio in the next study could fall anywhere in the range of 0.14 to 1.77. The value of  $I^2$  is 92.12, which means that 92% of the observed variance comes from real differences between studies and, as such, can potentially be explained by study-level covariates.

The next issue for the authors was to try and explain some of this variation. There was reason to believe that the drug was more effective in colder climates. This hypothesis was based on the theory that persons in colder climates were less likely to have a natural immunity to TB. It was also based on the expectation that the drug would be more potent in the colder climates (in warmer climates the heat would cause the drug to lose potency, and direct exposure to sunlight could kill some of the bacteria that were required for the vaccine to work properly).

In the absence of better predictor variables (such as the actual storage conditions used for the vaccine) Berkey *et al.* (1995) used absolute distance from the equator as a surrogate for climate (i.e. geographical regions in the Northern US would be colder than those in the tropics), and used regression to look for a relationship between *Distance* and treatment effect. Given the post hoc nature of this analysis, a positive finding would probably not be definitive, but would suggest a direction for additional research. (See also Sutton *et al.*, 2000; Egger *et al.*, Forthcoming.)

### Assessing the impact of the slope

Table 22.1 shows the data for each study (events and sample size, effect size and latitude). Table 22.2 shows the results for a meta-regression using absolute latitude to predict the log risk ratio.

**Table 22.1** The BCG dataset.

	Vaccinated		Control		RR	In(RR)	V <sub>InRR</sub>	Latitude
	TB	Total	TB	Total				
Vandiviere et al., 1973	8	2545	10	629	0.198	-1.621	0.223	19
Ferguson & Simes, 1949	6	306	29	303	0.205	-1.585	0.195	55
Hart & Sutherland, 1977	62	13598	248	12867	0.237	-1.442	0.020	52
Rosenthal et al., 1961	17	1716	65	1665	0.254	-1.371	0.073	42
Rosenthal et al., 1960	3	231	11	220	0.260	-1.348	0.415	42
Aronson, 1948	4	123	11	139	0.411	-0.889	0.326	44
Stein & Aaronson, 1953	180	1541	372	1451	0.456	-0.786	0.007	44
Coetzee & Berjak, 1968	29	7499	45	7277	0.625	-0.469	0.056	27
Comstock et al., 1974	186	50634	141	27338	0.712	-0.339	0.012	18
Frimodt-Møller et al., 1973.	33	5069	47	5808	0.804	-0.218	0.051	13
Comstock et al., 1976	27	16913	29	17854	0.983	-0.017	0.071	33
TB Prevention Trial, 1980	505	88391	499	88391	1.012	0.012	0.004	13
Comstock & Webster, 1969	5	2498	3	2341	1.562	0.446	0.533	33

**Table 22.2** Fixed-effect model – regression results for BCG.

Fixed effect, Z-distribution						
Point estimate	Standard error	95% Lower	95% upper	Z-value	p-Value	
Intercept	0.34356	0.08105	0.18471	0.50242	4.23899	0.00002
Latitude	-0.02924	0.00265	-0.03444	-0.02404	-11.02270	0.00000

The regression coefficient for latitude is -0.0292, which means that every one degree of latitude corresponds to a decrease of 0.0292 units in effect size. In this case, the effect size is the log risk ratio, and (given the specifics of this example) this corresponds to a more effective vaccination.

If we were working with a primary study and wanted to test the coefficient for significance we might use a *t*-test of the form

$$t = \frac{B}{SE_B}, \quad (22.1)$$

In meta-analysis the coefficient for any covariate (*B*) and its standard error are based on groups of *studies* rather than groups of *subjects* but the same logic applies. Historically,

**Table 22.3** Fixed-effect model – ANOVA table for BCG regression.

Analysis of variance			
	Q	df	p-value
Model ( <i>Q<sub>model</sub></i> )	121.49992	1	0.00000
Residual ( <i>Q<sub>resid</sub></i> )	30.73309	11	0.00121
Total ( <i>Q</i> )	152.23301	12	0.00000

in meta-regression the test is based on the Z-distribution, and that is the approach presented here (however, see notes at the end of this chapter about other approaches). The statistic to test the significance of the slope is

$$Z = \frac{B}{SE_B}. \quad (22.2)$$

Under the null hypothesis that the coefficient is zero,  $Z$  would follow the normal distribution.

In the running example the coefficient for latitude is  $-0.02924$  with standard error  $0.00265$ , so

$$Z = \frac{-0.02924}{0.00265} = -11.0227.$$

The two-tailed  $p$ -value corresponding to  $Z = -11.0227$  is  $< 0.00001$ . This tells us that the slope is probably not zero, and the vaccination is more effective when the study is conducted at a greater distance from the equator.

The Z-test can be used to test the statistical significance of any single coefficient but when we want to assess the impact of several covariates simultaneously we need to use the Q-test. (This is analogous to the situation in primary studies where we use a  $t$ -test to assess the impact of one coefficient but an  $F$ -test to assess the impact of two or more.)

As we did when working with analysis of variance we can divide the sum of squares into its component parts, and create an analysis of variance table as follows.

As before,  $Q$  is defined as a weighted sum of squares, which we can partition into its component parts.  $Q$  reflects the total dispersion of studies about the grand mean.  $Q_{\text{resid}}$  reflects the distance of studies from the regression line.  $Q_{\text{model}}$  is the dispersion explained by the covariates.

Each  $Q$  statistic is evaluated with respect to its degrees of freedom, as follows.

- $Q$  is 152.2330, with 12 degrees of freedom and  $p < 0.00001$  (this is the same value presented for the initial meta-analysis with no covariates). This means that the amount of total variance is more than we would expect based on within-study error.
- $Q_{\text{model}}$  is 121.4999 with 1 degree of freedom and  $p < 0.00001$ . This means that the relationship between latitude and treatment effect is stronger than we would expect by chance.
- $Q_{\text{resid}}$  is 30.7331 with 11 degrees of freedom and  $p < 0.0001$ . This means that even with latitude in the model, some of the between-studies variance (reflected by the distance between the regression line and the studies) remains unexplained.

$Q_{\text{model}}$  here is analogous to  $Q_{\text{bet}}$  for subgroup analysis, and  $Q_{\text{resid}}$  is analogous to  $Q_{\text{within}}$  for subgroup analysis. If the covariates are coded to represent subgroups, then  $Q_{\text{model}}$  will be identical to  $Q_{\text{bet}}$ , and  $Q_{\text{resid}}$  will be identical to  $Q_{\text{within}}$ .

### The Z-test and the Q-test

In this example there is only one covariate and so we have the option of using either the Z-test or the Q-test to test its relationship with effect size. It follows that the two

tests should yield the same results, and they do. The  $Z$ -value is  $-11.0227$ , with a corresponding  $p$ -value of  $<0.0001$ . The  $Q$ -value is  $121.4999$  with a corresponding  $p$ -value of  $<0.0001$  (with 1  $df$ ). Finally,  $Q$  should be equal to  $Z^2$  (since  $Q$  squares each difference while  $Z$  does not) and in fact  $-11.0227^2$  is equal to  $121.4999$ .

When we have more than one covariate the  $Q$  statistic serves as an omnibus test of the hypothesis that all the  $B$ s are zero. The  $Z$ -test can be used to test any coefficient, with the others held constant.

### Quantify the magnitude of the relationship

The  $Z$ -test, like all tests of significance, speaks to the question of statistical, rather than substantive, significance. Therefore, in addition to reporting the test of significance, one should always report the magnitude of the relationship. Here, the relationship of latitude to effect (expressed as a log risk ratio) is

$$\ln(RR) = 0.3435 - 0.0292(X)$$

where  $X$  is the absolute latitude. Figure 22.2 shows the plot of log risk ratio on latitude.

In the graph, each study is represented by a circle that shows the actual coordinates (observed effect size by latitude) for that study. The size (specifically, the area) of each circle is proportional to that study's weight in the analysis. Since this analysis is based on the fixed-effect model, the weight is simply the inverse of the within-study variance for each study.

The center line shows the predicted values. A study performed relatively close to the equator (such as the study performed in Madras, India, latitude 13) would have an expected effect near zero (corresponding to a risk ratio of 1.0, which means that the vaccination has no impact on TB). By contrast, a study at latitude 55 (Saskatchewan)

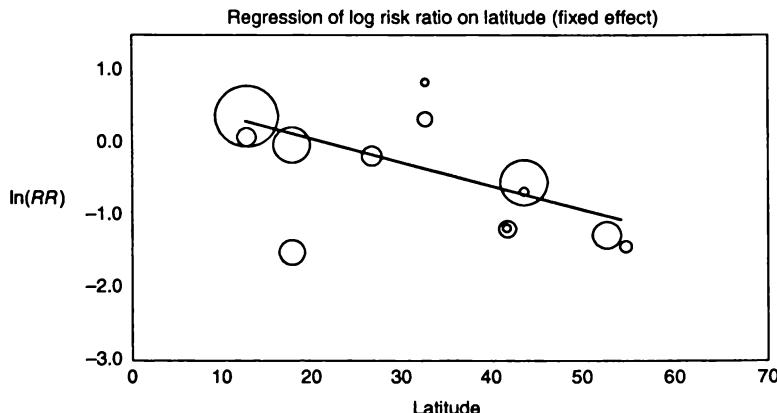


Figure 22.2 Fixed-effect model – regression of log risk ratio on latitude.

would have an expected effect near  $-1.20$  (corresponding to a risk ratio near  $0.30$ , which means that the vaccination is expected to decrease the risk of TB by about 70%).

The 95% confidence interval for  $B$  is given by

$$LL_B = B - 1.96 \times SE_B \quad (22.3)$$

and

$$UL_B = B + 1.96 \times SE_B. \quad (22.4)$$

In the BCG example

$$LL_B = (-0.0292) - 1.96 \times 0.0027 = -0.0344$$

and

$$UL_B = (-0.0292) + 1.96 \times 0.0027 = -0.0240.$$

In words, the true coefficient could be as low as  $-0.0344$  and as high as  $-0.0240$ . These limits can be used to generate confidence intervals on the plot.

## FIXED OR RANDOM EFFECTS FOR UNEXPLAINED HETEROGENEITY

In Part 3 we discussed the difference between a fixed-effect model and a random-effects model. Under the fixed-effect model we assume that the true effect is the same in all studies. By contrast, under the random-effects model we allow that the true effect may vary from one study to the next.

When we were working with a single population the difference in models translated into one true effect versus a distribution of true effects for all studies. When we were working with subgroups (analysis of variance) it translated into one true effect versus a distribution of effects for all studies within a subgroup (for example, studies that used Intervention A or studies that used Intervention B). For meta-regression it translates into one true effect versus a distribution of effects for studies with the same predicted value (for example, two studies that share the same value on all covariates). This is shown schematically in Figure 22.3 and Figure 22.4.

Under the fixed-effect model, for any set of covariate values (that is, for any predicted value) there is one population effect size (represented by a circle in Figure 22.3). Under the random-effects model, for any predicted value there is a distribution of effect sizes (in Figure 22.4 the distribution is centered over the predicted value but the population effect size can fall to the left or right of center). (In both figures we assume that the prediction is perfect, so that the true effect (or mean effect) falls directly on the prediction line.)

As always, the selection of a method should follow the logic of how the studies were selected. When we introduced the idea of fixed versus random effects for a single population, the example we used for the fixed-effect model was a pharmaceutical company that planned a series of ten trials that were identical for all intents and purposes (page 76). When we moved on to analysis of variance we extended the same example, and assumed that five cohorts would be used to test placebo versus *Drug A*, while five would be used to test placebo versus *Drug B* (page 172). We can extend this example, and use it to create an example where a fixed-effect analysis would be

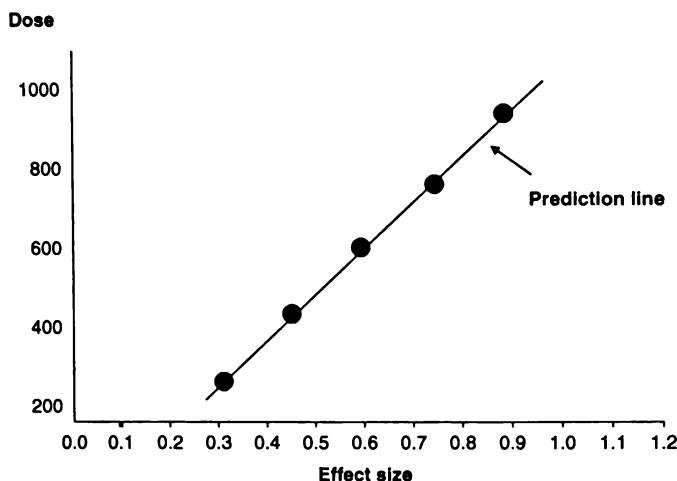


Figure 22.3 Fixed-effect model – population effects as function of covariate.

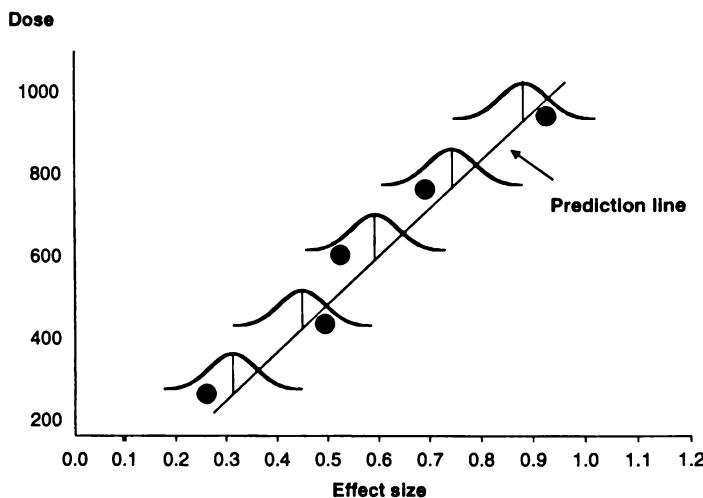


Figure 22.4 Random-effects model – population effects as a function of covariate.

appropriate for meta-regression. As before, we'll assume a total of 10 studies, five for placebo versus *Drug A* and five for placebo versus *Drug B*. This time we'll assume that each study used either 10, 20, 40, 80, or 160 mg. of the drug and we will use regression to assess the impact of drug and dose. The fixed-effect model makes sense here because the studies are known to be identical on all other factors.

As before, we note that this example is not representative of most systematic reviews. In the vast majority of cases, especially when the studies are performed by different researchers and then culled from the literature, it is more plausible that the impact of the covariates captures some, *but not all*, of the true variation among effects. In this case, it is the random-effects model that reflects the nature of the distribution of true effects, and should therefore be used in the analyses. Also as before, if the study design suggests that the random-effects model is appropriate, then this model should be selected *a priori*. It is a mistake to start with the fixed-effect model and move on to random effects only if the test for heterogeneity is statistically significant.

Since the meaning of a summary effect size is different for fixed versus random effects, the null hypothesis being tested also differs. Both test a null hypothesis of no linear relationship between the covariates and the effect size. The difference is that under the fixed-effect model that effect size is the common effect size for all studies with a given value of the covariates. Under the random-effects model that effect size is the mean of the true effect sizes for all studies with a given value of the covariates. This is important, because the different null hypotheses reflect different assumptions about the sources of error. This means that different error terms are used to compute tests of significance and confidence intervals under the two models.

Computationally, the difference between fixed effect and random effects is in the definition of the variance. Under the fixed-effect model the variance is the variance within studies, while under the random-effects model it is the variance within studies plus the variance between studies ( $\tau^2$ ). This holds true for one population, for multiple subgroups, and for regression, but the mechanism used to estimate  $\tau^2$  depends on the context. When we are working with a single population,  $\tau^2$  reflects the dispersion of true effects across all studies, and is therefore computed for the full set of studies. When we are working with subgroups,  $\tau^2$  reflects the dispersion of true effects within a subgroup, and is therefore computed within subgroups. When we are working with regression,  $\tau^2$  reflects the dispersion of true effects for studies with the same predicted value (that is, the same value on the covariates) and is therefore computed for each point on the prediction slope. As a practical matter, of course, most points on the slope have only a single study, and so this computation is less transparent than that for the single population (or subgroups) but the concept is the same. The computational details are handled by software and will not be addressed here.

The practical implications of using a random-effects model rather than a fixed-effect model for regression are similar to those that applied to a single population and to subgroups. First, the random-effects model will lead to more moderate weights being assigned to each study. As compared with a fixed-effect model, the random-effects model will assign more weight to small studies and less weight to large studies. Second, the confidence interval about each coefficient (and slope) will be wider than it would be under the fixed-effect model. Third, the *p*-values corresponding to each coefficient and to the model as a whole are less likely to meet the criterion for statistical significance.

As always, the selection of a model must be based on the context and characteristics of the studies. In particular, if there is heterogeneity in true effects that is not explained by the covariates, then the random-effects model is likely to be more appropriate.

### RANDOM-EFFECTS MODEL

We return now to the BCG example and apply the random-effects model (Figure 22.5).

Using a random effects analysis the risk ratio for the 13 studies is 0.490 with a confidence interval of 0.345 to 0.695, which says that the *mean effect* of the vaccine was to decrease the risk of TB by at least 30% and possibly by as much as 65% (see Figure 22.5). In log units the risk ratio is -0.714 with a standard error of 0.179, which yields a Z-value of -3.995 ( $p < 0.001$ ) which allows us to reject the null hypothesis of no effect.

At least as interesting, however, is the substantial variation in the treatment effects, which ranged from a risk ratio of 0.20 (vaccine *reduces* the risk by 80%) to 1.56 (vaccine *increases* the risk by 56%). The relevant statistics were presented when we discussed the fixed-effect model.

### Assessing the impact of the slope

A meta-regression using the random-effect model (method of moments) yields the results shown in Table 22.4.

This has the same format as Table 22.2, showing the coefficients for predicting the log risk ratio from latitude and related statistics. We use an asterisk (\*) to indicate

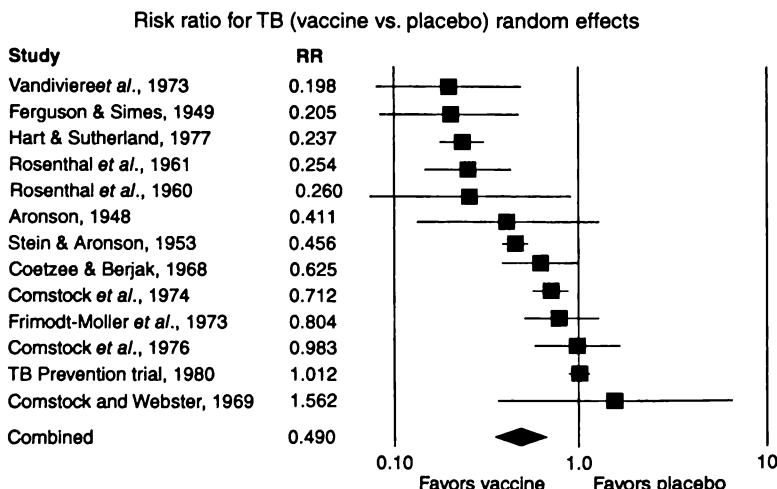


Figure 22.5 Random-effects model – forest plot for the BCG data.

**Table 22.4** Random-effects model – regression results for BCG.

Random effects, Z-distribution						
	Point estimate	Standard error	95% lower	95% upper	Z-value	p-Value
Intercept	0.25954	0.23231	-0.19577	0.71486	1.11724	0.26389
Latitude	-0.02923	0.00673	-0.04243	-0.01603	-4.34111	0.00001

that these statistics are based on the random-effects model. With that distinction, the interpretation of the slope(s) is the same as that for the fixed-effect model. Concretely,

$$Z^* = \frac{B^*}{SE_{B^*}}. \quad (22.5)$$

Under the null hypothesis that the coefficient is zero,  $Z^*$  would follow the normal distribution.

In the running example the coefficient for latitude is -0.0923 with standard error 0.00673, so

$$Z = \frac{-0.02923}{0.00673} = -4.3411.$$

(In this example the slope happens to be almost identical under the fixed-effect and random-effects models, but this is not usually the case.) The two-tailed  $p$ -value corresponding to  $Z^* = -4.3411$  is 0.00001. This tells us that the slope is probably not zero, and the vaccination is more effective when the study is conducted at a greater distance from the equator.

Under the null hypothesis that none of the covariates 1 to  $p$  is related to effect size,  $Q_{model}^*$  would be distributed as chi-squared with degrees of freedom equal to  $p$ . In the running example,  $Q_{model}^* = 18.8452$ ,  $df = 1$ , and  $p = 0.00001$  (see Table 22.5).

In this example there is only one covariate (latitude) and so we have the option of using either the  $Z$ -test or the  $Q$ -test to assess the impact of this covariate. It follows that the two tests should yield the same results, and they do. The  $Z$ -value is -4.3411, with a corresponding  $p$ -value of 0.00001. The  $Q$ -value is 18.8452 with a corresponding  $p$ -value of 0.00001. Finally,  $Q_{model}^*$  should be equal to  $Z^2$  and in fact 18.8452 equals -4.3411<sup>2</sup>.

The goodness of fit test addresses the question of whether there is heterogeneity that is not explained by the covariates.  $Q_{resid}$  can also be used to estimate (and test) the

**Table 22.5** Random-effects model – test of the model.

#### Test of the model:

**Simultaneous test that all coefficients (excluding intercept) are zero**

$Q_{model}^* = 18.8452$ ,  $df = 1$ ,  $p = 0.00001$

**Goodness of fit: test that unexplained variance is zero**

$T^2 = 0.063$ ,  $SE = 0.055$ ,  $Q_{resid} = 30.733$ ,  $df = 11$ ,  $p = 0.00121$

variance,  $\tau^2$ , of this unexplained heterogeneity. This  $Q_{\text{resid}}$  is the weighted residual SS from the regression using fixed-effect weights (see Table 22.3)

$Q_{\text{model}}^*$  here is analogous to  $Q_{\text{het}}^*$  for subgroup analysis, and  $Q_{\text{resid}}$  is analogous to  $Q_{\text{within}}$  for subgroup analysis. If the covariates represent subgroups, then  $Q_{\text{model}}^*$  is identical to  $Q_{\text{het}}^*$  and  $Q_{\text{resid}}$  is identical to  $Q_{\text{within}}$ . If there are no predictors then  $Q_{\text{resid}}$  here is the same as  $Q$  for the original meta-analysis.

When working with meta-regression with the fixed-effect model we were able to partition the total variance into a series of components, with  $Q_{\text{model}}$  plus  $Q_{\text{resid}}$  summing to  $Q$ . This was possible with the fixed-effect model because the weight assigned to each study was determined solely by the within-study error, and was therefore the same for all three sets of calculations. By contrast, under the random effects model the weight assigned to each study incorporates between-studies variance also, and this varies from one set of calculations to the next. Therefore, the variance components are not additive. For that reason, we display an analysis of variance table for the fixed-effect analysis, but not for the random-effects analysis.

### Quantify the magnitude of the relationship

The relationship of latitude to effect (expressed as a log risk ratio) is

$$\ln(RR) = 0.2595 - 0.0292(X)$$

where X is the absolute latitude. We can plot this in Figure 22.6.

In this Figure, each study is represented by a circle that shows the actual coordinates (observed effect size by latitude) for that study. The size (specifically, the area) of each circle is proportional to that study's weight in the analysis. Since this analysis is based on the random-effects model, the weight is the total variance (within-study plus  $\tau^2$ ) for each study.

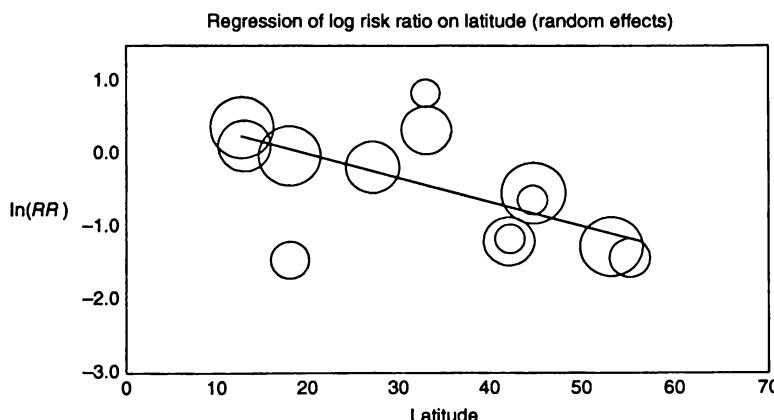


Figure 22.6 Random-effects model – regression of log risk ratio on latitude.

Note the difference from the fixed-effect graph (Figure 22.2). When using random effects, the weights assigned to each study are more similar to one another. For example, the TB prevention trial (1980) study dominated the graph under the fixed-effect model (and exerted substantial influence on the slope) while Comstock and Webster (1969) had only a trivial impact (the relative weights for the two studies are 41% and 0.3% respectively). Under random effects the two are more similar (14% and 1.6% respectively).

The center line shows the predicted values. A study performed relatively close to the equator (latitude of 10) would have an expected effect near zero (corresponding to a risk ratio of 1.0, which means that the vaccination has no impact on TB). By contrast, a study at latitude 55 (Saskatchewan) would have an expected effect near -1.50 (corresponding to a risk ratio near 0.20, which means that the vaccination decreased the risk of TB by about 80%).

The 95% confidence interval for  $B$  is given by

$$LL_B = B^* - 1.96 \times SE_B \quad (22.6)$$

and

$$UL_B = B^* + 1.96 \times SE_B. \quad (22.7)$$

In the running example

$$LL_B = (-0.0292) - 1.96 \times 0.0067 = -0.0424$$

and

$$UL_B = (-0.0292) + 1.96 \times 0.0067 = -0.0160.$$

In words, the true coefficient could be as low as -0.0424 and as high as -0.0160.

### The proportion of variance explained

In Chapter 19 we introduced the notion of the proportion of variance explained by subgroup membership in a random-effects analysis. The same approach can be applied to meta-regression.

Consider Figure 22.7, which shows a set of six studies with no covariate. Since there is no covariate the prediction slope is simply the mean (the intercept, if we were to compute a regression), depicted by a vertical line. The normal distribution at the bottom of the figure reflects  $T^2$ , and is needed to explain why the dispersion *from the prediction line* (the mean) exceeds the within-study error.

Now, consider Figure 22.8, which shows the same size studies with a covariate  $X$ , and the prediction slope depicted by a line that reflects the prediction equation. The normal distribution at each point on the prediction line reflects the value of  $T^2$ , and is needed to explain why the dispersion *from the prediction line* (this time, the slope) exceeds the within study error. Because the covariate explains some of the between-studies variance, the  $T^2$  in Figure 22.8 is smaller than the one in Figure 22.7, and the ratio of the two can be used to quantify the proportion of variance explained.

Note 1. Normally, we would plot the effect size on the  $y$ -axis and the covariate on the  $x$ -axis (see, for example, Figure 22.6). Here, we have transposed the axes to maintain the parallel with the forest plot.

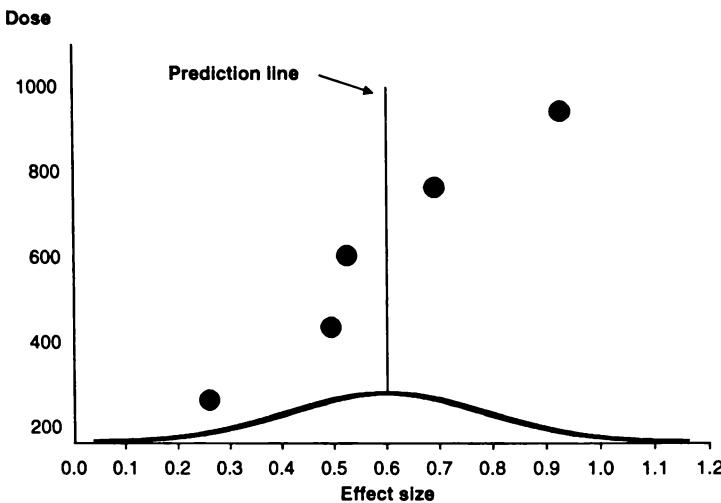


Figure 22.7 Between-studies variance ( $T^2$ ) with no covariate.

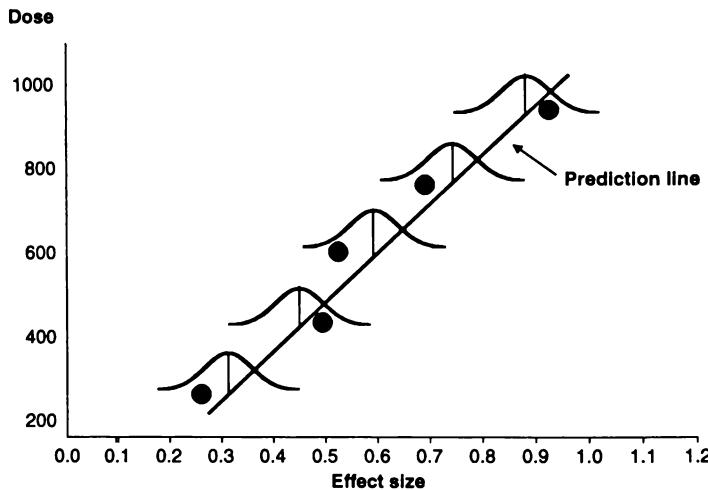


Figure 22.8 Between-studies variance ( $T^2$ ) with covariate.

Note 2. For clarity, we plotted the true effects for each figure. In practice, of course, we observe estimates of the true effects, remove the portion of variance attributed to within-study error, and impute the amount of variance remaining.

In primary studies, a common approach to describing the impact of a covariate is to report the proportion of variance explained by that covariate. That index,  $R^2$ , is defined

as the ratio of explained variance to total variance,

$$R^2 = \frac{\sigma_{explained}^2}{\sigma_{total}^2} \quad (22.8)$$

or, equivalently,

$$R^2 = 1 - \left( \frac{\sigma_{unexplained}^2}{\sigma_{total}^2} \right) \quad (22.9)$$

This index is intuitive as it can be interpreted as a ratio, with a range of 0 to 1, or of 0% to 100%. Many researchers are familiar with this index, and have a sense of what proportion of variance is likely to be explained by different kinds of covariates or interventions.

This index cannot be applied directly to meta-analysis. The reason is that in meta-analysis the total variance includes both variance within studies and between studies. The covariates are study-level covariates, and as such they can potentially explain only the between-studies portion of the variance. In the running illustration, even if the *true* effect for each study fell directly on the prediction line the proportion of variance explained would not approach 1.0 because the *observed* effects would fall at some distance from the prediction line.

Therefore, rather than working with this same index we use an analogous index, defined as the *true* variance explained, as a proportion of the total *true* variance. Since the true variance is the between-studies variance,  $\tau^2$ , we compute

$$R^2 = \frac{T_{explained}^2}{T_{total}^2} \quad (22.10)$$

or

$$R^2 = 1 - \left( \frac{T_{unexplained}^2}{T_{total}^2} \right). \quad (22.11)$$

In the running example (Table 22.6)  $T_{total}^2$  for the full set of studies was 0.309, and  $T_{unexplained}^2$  for the equation with latitude is 0.063. This gives us

$$R^2 = 1 - \left( \frac{0.063}{0.309} \right) = 0.7950 \quad (22.12)$$

**Table 22.6** Random-effects model – comparison of model (latitude) versus the null model.

**Comparison of model with latitude versus the null model**

Total between-study variance (intercept only)

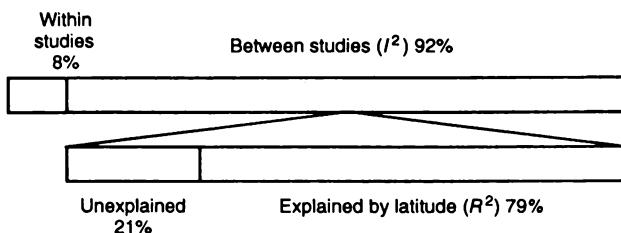
$T_{total}^2 = 0.309$ , SE = 0.230,  $Q_{resid} = 152.233$ , df = 12,  $p = 0.00000$

Unexplained between-study variance (with latitude in model)

$T_{unexplained}^2 = 0.063$ , SE = 0.055,  $Q_{resid} = 30.733$ , df = 11,  $p = 0.0012$

Proportion of total between-study variance explained by the model

$R^2$  analog = 1 - (0.063/0.309) 79.50%



**Figure 22.9** Proportion of variance explained by latitude.

This is shown schematically for the running example (see Figure 22.9). In Figure 22.9, the top line shows that 8% of the total variance was within studies and 92% was between studies (which is also the definition of  $I^2$ ). The within-studies variance cannot be explained by a study-level moderator, and so is removed from the equation and we focus on the shaded part.

On the bottom line, the type of intervention is able to explain 79% of the relevant variance, leaving 21% unexplained. Critically, the 79% and 21% sum to 100%, since we are concerned only with the variance between-studies.

While the  $R^2$  index has a range of 0 to 1 in the population, it is possible for sampling error to yield a value of  $R^2$  that falls outside of this range. In that case, the value is set to either 0 (if the estimate falls below 0) or to 1 (if it falls above 1).

### SUMMARY POINTS

- Just as we can use multiple regression in primary studies to assess the relationship between subject-level covariates and an outcome, we can use meta-regression in meta-analysis to assess the relationship between study level covariates and effect size.
- Meta-regression may be performed under the fixed-effect or the random effects model, but in most cases the latter is appropriate.
- In addition to testing the impact of covariates for statistical significance, it is important to quantify the magnitude of their relationship with effect size. For this purpose we can use an index based on the percent reduction in true variance, analogous to the  $R^2$  index used with primary studies.

# Notes on Subgroup Analyses and Meta-Regression

---

Introduction

Computational model

Multiple comparisons

Software

Analyses of subgroups and regression analyses are observational

Statistical power for subgroup analyses and meta-regression

---

## INTRODUCTION

In this chapter we address a number of issues that are relevant to both subgroup analyses (analysis of variance) and to meta-regression.

## COMPUTATIONAL MODEL

The researcher must always choose between a fixed-effect model and a random-effects model. When we are working with a single set of studies the fixed-effect analysis assumes that all studies share a common effect size. When we are working with subgroups, it assumes that all studies within a subgroup share a common effect size. When we are working with meta-regression, it assumes that all studies which have the same values on the covariates share a common effect size.

These kinds of assumption can sometimes be justified, as in the pharmaceutical example that we used on pages 76, 172, and 203. In most cases, however, especially when the studies for the review have been culled from the literature, it is more plausible to assume that the subgroup membership or the covariates explain *some*, but *not all*, of the dispersion in effect sizes. Therefore, the random-effects model is more likely to fit the data, and is the model that should be selected.

### Mistakes to avoid in selecting a model

When we introduced the idea of a random-effects model in Chapter 12 we noted that researchers sometimes start with a fixed-effect model and then move to a random effects model if there is empirical evidence of heterogeneity (a statistically significant  $p$ -value). In the case of subgroup analysis this approach would suggest that we start by using the fixed-effect model within groups, and then move to the random-effects model only if  $Q$  within groups was statistically significant. In the case of meta-regression it would suggest that we start by using the fixed-effect model, and then move to the random-effects model only if the  $Q$  for the residual error was statistically significant.

We explained that this approach was problematic when working with a single set of studies, and it continues to be a bad idea here, for the same reasons. If substantive considerations suggest that the effect size is likely to vary (within the full set of studies, within subgroups, or for studies with a common set of covariate values) then we should be using the corresponding model even if the test for heterogeneity fails to yield a significant  $p$ -value. This lack of significance means only that we have failed to meet a certain threshold of proof (possibly because of low statistical power) and does not prove that the studies share a common effect size.

### Practical differences related to the model

Researchers often ask about the practical implications of using a random-effects model rather than a fixed-effect model. The random-effects model will apportion the study weights more evenly, so that a large study has less impact on the summary effect (or regression line) than it would have under the fixed-effect model, and a small study has more impact than it would have under the fixed-effect model. Also, confidence intervals will tend to be wider under the random-effects model than under the fixed-effect model. While this tells us what the impact will be of using the fixed-effect or random-effects model, it says nothing about which model we should use. The only issue relevant to that decision is the question of which model fits the data.

### The null hypothesis under the different models

Since the meaning of a summary effect size is different for fixed versus random effects, the null hypothesis being tested also differs under the two models.

Recall that when we are working with a single group, under the fixed-effect model the summary effect size represents the common effect size for all the studies. The null hypothesis is that the common effect size is equal to a nil value (0.00 for a difference, or 1.00 for a ratio). By contrast, under the random-effects model the summary effect size represents the mean of the distribution of effect sizes. The null hypothesis is that the mean of all possible studies is equal to the nil value.

In subgroup analyses, under the fixed-effect model the summary effect size for subgroups  $A$  and  $B$  each represents the common effect size for a group of studies. The null hypothesis is that the common effect for the  $A$  studies is the same as the common

effect for the *B* studies. By contrast, under the random-effects model the effect size for subgroups *A* and *B* each represents the mean of a distribution of effect sizes. The null hypothesis is that the mean of all possible *A* studies is identical to the mean of all possible *B* studies.

In regression, under the fixed-effect model we assume that there is one true effect size for any given value of the covariate(s). The null hypothesis is that this effect size is the same for all values of the covariate(s). By contrast, under the random effects model we assume that there is a distribution of effect sizes for any given value of the covariate(s). The null hypothesis is that this mean is the same for all values of the covariate(s).

While the distinction between a *common* effect size and a *mean* effect size might sound like a semantic nuance, it actually reflects an important distinction between the models. In the case of the fixed-effect model, because we assume that we are dealing with a common effect size, we apply an error model which assumes that the between-studies variance is zero. In the case of the random-effects model, because we allow that the effect sizes may vary, we apply an error model which makes allowance for this additional source of uncertainty. This difference has an impact on the mean (the summary effect for a single group, the summary effect within subgroups, and the slope in a meta-regression). It also has an impact on the standard error, tests of significance and confidence intervals.

### Some technical considerations in random-effects meta-regression

As is the case for a standard random-effects meta-analysis in the absence of covariates, several methods are available for estimating  $\tau^2$  in meta-regression, including a moment method and maximum likelihood method. In practice, any differences among methods will usually be small.

The results we presented in this chapter used a moment estimate, which is the same as the method we used in Chapter 16 to estimate  $\tau^2$  for a single group. If we were to perform a meta-regression with no covariates, our estimate of  $\tau^2$  would be the same as the estimate we would obtain using the formulas in that chapter.

Whichever method is used to estimate  $\tau^2$ , the use of a Z-test to assess the statistical significance of a covariate (or the difference between two subgroups), while common, is not strictly appropriate. When dealing with simple numerical data, to compute a confidence interval or to test the significance of a difference (or variance) we use *Z* if the sampling distribution is *known*. By contrast, we use *t* if the sampling distribution is being *estimated* from the dispersion observed in the sample (as it is, for example, when we compare means using a *t*-test).

Similarly, in meta-analysis, the *Z*-distribution is appropriate only for the fixed-effect model, where the only source of error is within studies. By contrast, when we use a random-effects model, we are estimating the dispersion across studies, and should account for this by using a *t*-distribution. Several methods have been proposed to address this issue, including one by Knapp and Hartung (2003) which is outlined below. While these can be applied to any use of the random-effects model (for a single group of studies, for a subgroup analysis, and for meta-regression), they have to date only been implemented in software for meta-regression.

The Knapp–Hartung method involves two modifications to the standard error for the random-effects model. First, the between-studies component of the variance is multiplied by a factor that makes it correspond to the  $t$ -distribution rather than the  $Z$ -distribution. Second, the test statistic is compared against the  $t$ -distribution rather than the  $Z$ -distribution. This has the effect of expanding the width of the confidence intervals and of moving the  $p$ -value away from zero.

Higgins and Thompson (2004) proposed an approach that bypasses the sampling distributions and instead employs a permutation test to yield a  $p$ -value. Using this approach we compute the  $Z$ -score corresponding to the observed covariate. Then, we randomly redistribute the covariates among studies and see what proportion of these yield a  $Z$ -score exceeding the one that we had obtained. This proportion may be viewed as an exact  $p$ -value.

## MULTIPLE COMPARISONS

In primary studies researchers often need to address the issue of multiple comparisons. The basic problem is that if we conduct a series of analyses with alpha set at 0.05 for each, then the overall likelihood of a type I error (assuming that the null hypothesis is actually true) will exceed 5%. This problem crops up when a study includes more than two groups and we compare more than one pair of means. It also arises when we perform an analysis on more than one outcome.

While there is consensus that conducting many comparisons can pose a problem, there is no consensus about how this problem should be handled. Some suggest conducting an omnibus test that asks if there are any nonzero effects, and then proceeding to look at pair-wise comparisons only if the initial test meets the criterion for significance. Others suggest going straight to the pairwise tests but using a stricter criterion for significance (say 0.01 rather than 0.05 for five tests). Hedges and Olkin (1985) discuss this and other methods to control the error rate when using multiple tests. Some suggest that the researcher not make any formal adjustment, but evaluate the data in context. For example, one significant  $p$ -value in forty tests would not be seen as grounds for rejection of the null hypothesis.

Essentially the same issue exists in meta-analysis. In the case of subgroup analyses, if a meta-analysis includes a number of subgroups, the issue of multiple comparisons arises when we start to compare several pairs of subgroup means. In the case of meta-regression this issue arises when we include a number of covariates and want to test each one for significance. As with primary studies, while there is consensus that conducting many comparisons can pose a problem, there is no consensus about how this problem should be handled. The approaches generally used for primary studies can be applied to meta-analysis as well.

## SOFTWARE

Some of the programs developed for meta-analysis are able to perform subgroup analysis as well as meta-regression (see Chapters 49 and 51). Note that programs intended

for statistical analysis of primary studies should not be used to perform these procedures in meta-analysis, for two reasons. First, routines for analysis of variance or multiple regression intended for primary studies do not weight the studies, as is needed for meta-analysis. While most programs do allow the user to assign weights, this becomes a difficult procedure when we move to random-effects weights (which are usually the ones we want to use). Second, the rules for assigning degrees of freedom in the analysis of variance and meta-regression are different for meta-analysis than for primary studies, and so using the primary-study routines for a meta-analysis will yield incorrect standard errors and *p*-values.

## ANALYSES OF SUBGROUPS AND REGRESSION ANALYSES ARE OBSERVATIONAL

In a randomized trial, participants are assigned at random to a condition (such as treatment versus placebo). Because the participants are assumed to be similar in all respects except for the treatment condition, differences that do emerge between conditions can be attributed to the treatment. By contrast, in an observational study we compare pre-existing groups, such as workers with a college education versus those who did not attend college. While we can report on differences in wages of the two groups we cannot attribute this outcome to the amount of schooling because the groups differ in various ways. For example, we are likely to find that those with a college education are paid more, but we cannot attribute this to their schooling since it could be due (at least in part) to other factors associated with higher socioeconomic status.

The issue of randomized versus observational studies as it relates to meta-analysis is discussed in Chapter 45. There, we discuss the fact that randomized studies and observational studies address different questions, and for this reason it generally makes sense to include only one or the other in a given meta-analysis.

However, there is one issue that is directly relevant to the present discussion, as follows. Assume we start with a set of randomized experiments that assess the impact of an intervention. The effect in any single experiment could serve to establish causality and the summary effect can also serve to establish causality. This is because the relationship between treatment and outcome is protected by the randomization process (it *must be* due to treatment) in each study, and this protection carries over to the summary effect.

However, *even if the individual studies are randomized trials*, once we move beyond the goal of reporting a summary effect and proceed to perform a subgroup analyses or meta-regression, *we have moved out of the domain of randomized experiments, and into the domain of observational studies*. For example, suppose that half the studies used a low dose of aspirin while half used a high dose, and that the impact of the treatment was significantly stronger in the high-dose studies. It is *possible* that the difference is due to the dose, but it is also possible that the studies that used a higher dose differed in some systematic way from the other studies. Perhaps these studies used patients who were in poor health, or older, and therefore more likely to benefit from the treatment. Therefore, the difference between subgroups, or the relationship between a

covariate and effect size, is *observational*. The same caveats that apply to any observational studies, in particular the fact that relationship does not imply causality, apply here too.

That said, in primary observational studies, researchers sometimes use regression analysis to try and remove the impact of potential confounders. In the aspirin example they might enter covariates in the sequence of health, age, and dose, to assess the impact of dose with health and age held constant. This is not a perfect solution since there may be other confounders of which we are not aware, but this approach can help to isolate the impact of specific factors and generate hypotheses to be tested in randomized trials. The same holds true for meta-regression. Of course, since covariate values are assigned at the study level, meta-regression can be used to adjust for potential confounders only for comparisons across studies, and not for potential confounders *within* studies.

There is one exception to the rule that subgroup analysis and regression cannot prove causality. This exception is the case where we know that the studies are identical in all respects except for the one captured by subgroup membership or by the covariate. The pharmaceutical example is a case in point. Here, we enrolled 1000 patients and assigned some to studies that would test a low dose of the drug vs. placebo, and others to studies that would test a high dose of the drug vs. placebo. Here, the assignment to subgroups is random. The same would apply if the patients were assigned to ten studies where the dose of drug was varied on a continuous scale, and we used meta-regression to test the relationship between dose and effect size. This set of circumstances is rarely (if ever) found in practice.

## STATISTICAL POWER FOR SUBGROUP ANALYSES AND META-REGRESSION

Statistical power is the likelihood that a test of significance will reject the null hypothesis. In the case of subgroup analyses it is the likelihood that the Z-test to compare the effect in two groups, or the Q-test to compare the effects across a series of groups, will yield a statistically significant  $p$ -value. In the case of meta-regression it is the likelihood that the Z-test of a single covariate or the Q-test of a set of covariates will yield a statistically significant  $p$ -value.

Power depends on the size of the effect and the precision with which we measure the effect. For subgroup analysis this means that power will increase as the difference between (or among) subgroup means increases, and/or the standard error within subgroups decreases. For meta-regression this means that power will increase as the magnitude of the relationship between the covariate and effect size increases, and/or the precision of the estimate increases. In both cases, a key factor driving the precision of the estimate will be the total number of individual subjects across all studies and (for random effects) the total number of studies.

While there is a general perception that power for testing the main effect is consistently high in meta-analysis, this perception is not correct (see Chapter 34) and certainly does not extend to tests of subgroup differences or to meta-regression. The

failure to find a statistically significant  $p$ -value when comparing subgroups or in meta-regression could mean that the effect (if any) is quite small, but could also mean that the analysis had poor power to detect even a large effect. One should never use a non-significant finding to conclude that the true means in subgroups are the same, or that a covariate is not related to effect size.

### SUMMARY POINTS

- The selection of a computational model (fixed-effect or random-effects) should be based on our understanding of the underlying distribution. In most cases, especially when the studies have been gathered from the published literature, the random-effects model (within-subgroups) is more plausible than the fixed-effect model.
- The strategy of starting with the fixed-effect model and then moving to the random-effects (or mixed-effect) model if the test for heterogeneity is significant, is a mistake, and should be strongly discouraged.
- The problem of performing multiple tests (the fear that the actual alpha may exceed the nominal alpha) is similar in meta-analysis to the same problem in primary studies, and similar strategies are suggested for dealing with this problem.
- The relationship between effect size and subgroup membership, or between effect size and covariates, is observational, and cannot be used to prove causality. This holds true even if all studies in the analysis are randomized trials. The protection afforded by the study design carries over to the summary effect across all studies, but not to other analyses.
- Statistical power for detecting a difference among subgroups, or for detecting the relationship between a covariate and effect size, is often low, and the usual caveats apply. To wit, failure to obtain a statistically significant difference among subgroups should never be interpreted as evidence that the effect is the same across subgroups. Similarly, failure to obtain a statistically significant effect for a covariate should never be interpreted as evidence that there is no relationship between the covariate and the effect size.

### Further Reading

- Borenstein, M. (2019). *Common Mistakes in Meta-Analysis and How to Avoid Them*. Englewood, NJ: Biostat, Inc.
- Cohen, J., West, S.G., Cohen, P., & Aiken, L. (2002). *Applied Multiple Regression/Correlation Analysis for the Behavioral Sciences* (3<sup>rd</sup> ed.). Mahwah, NJ, Lawrence Erlbaum Assoc.
- Higgins, J.P.T., & Thompson, S.G (2004). Controlling the risk of spurious findings from meta-regression. *Statistics in Medicine* 23, 1663–1682.
- Knapp, G. & Hartung, J. (2003). Improved tests for a random effects meta-regression with a single covariate. *Statistics in Medicine* 22, 2693–2710.



# Putting it all in Context



# Looking at the Whole Picture

---

### Introduction

Methylphenidate for adults with ADHD

Impact of GLP-1 mimetics on blood pressure

Augmenting clozapine with a second antipsychotic

Conclusions

Caveats

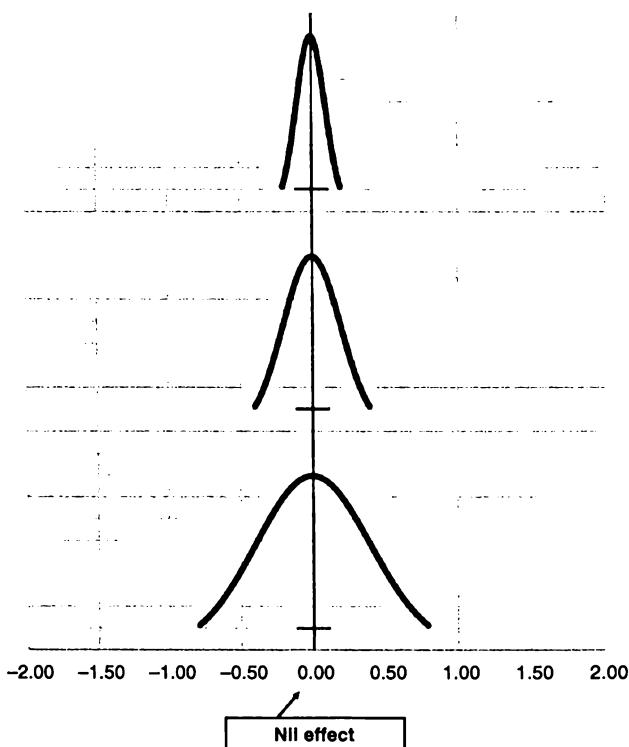
---

## INTRODUCTION

In a primary study, we tend to focus on the mean effect size and say nothing about how the effect size varies from study to study. This is reasonable, given that the primary study provides information about one study only. Unfortunately, many researchers adopt the same approach for a meta-analysis. They focus on the mean effect size. They may pay lip service to heterogeneity, remarking that it is *low* or *high*, but fail to quantify the heterogeneity in a meaningful way, nor to consider the practical implications of the heterogeneity. A key advantage of a systematic review and meta-analysis is that it allows us to understand how the effect size varies from study to study, and it is imperative that we take advantage of this opportunity.

One reason that many researchers fail to consider the practical implications of heterogeneity, is that they do not entirely understand heterogeneity. If a study reports that the mean effect size is a risk ratio of 0.70 and that  $I^2$  is 50%, the researcher might have no idea what the dispersion actually looks like. In that case, it follows that they cannot consider the implications of that dispersion. The first step in addressing this issue is to provide a mechanism for reporting dispersion in a way that the researcher and reader will understand. That mechanism is the prediction interval, as discussed in prior chapters. With that mechanism available to us, we can integrate information about the mean effect size and the dispersion of effects.

For example, consider Figures 24.1–24.3. Each figure displays the results of three fictional meta-analyses. In all cases, the effect size index is the standardized mean difference  $d$ . For purposes of this example, assume the following. An effect size of



**Figure 24.1** Three fictional examples where the mean effect is 0.00.

0.00 is a nil effect, which means that there is no difference between group means. An effect size of 0.20 is small, in the sense that it has little clinical or substantive importance. An effect size of 0.40 is moderate, in the sense that it has a moderate level of clinical or substantive importance. An effect size of 0.80 is large, in the sense that it has a high level of clinical or substantive importance.

The distributions differ from each other on two discrete dimensions. The first is the mean effect size, and the second is the dispersion of effects about the mean.

- The mean effect size varies by figure. Specifically, as we move from page to page, the mean effect size shifts from 0.0 to 0.40 to 0.80.
- The distribution of effects varies by panel within each figure. Specifically, as we move from (a) the top panel to (b) the middle panel to (c) the bottom panel, the effects (a) are consistent, (b) vary somewhat, or (c) vary substantially, across populations. To understand the potential utility of an intervention, we need to consider both the mean effect size and the distribution of effects.

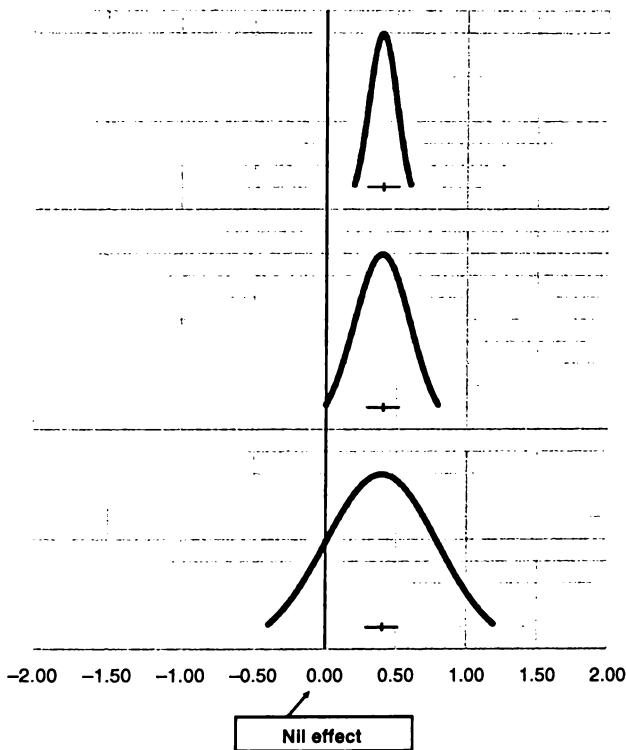
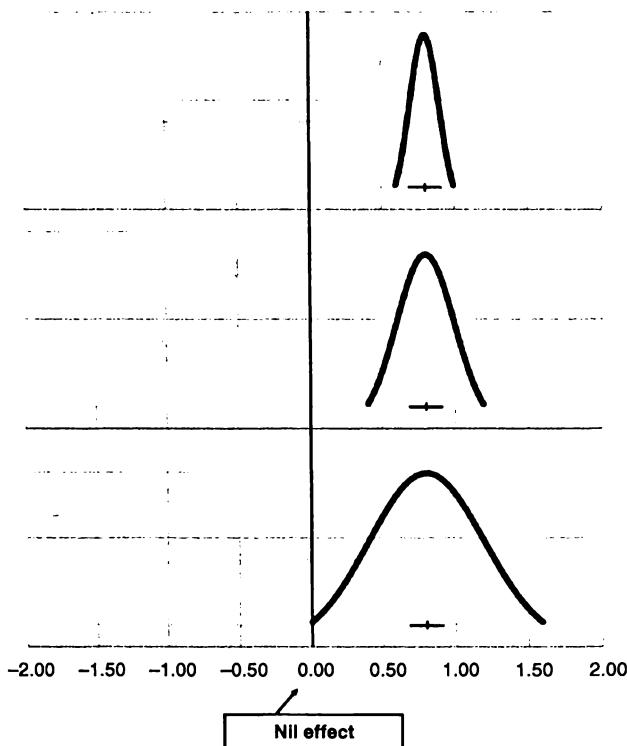


Figure 24.2 Three fictional examples where the mean effect is 0.40.

In Figure 24.1, the mean is always 0.00 and not statistically significant, but it is also important to distinguish among the three panels. In the top panel, the treatment is never effective. In the middle panel, there are some populations where the treatment is moderately harmful, and some where the treatment is moderately helpful. In the bottom panel, there are some populations where the treatment will cause major harm and some where it will be exceptionally helpful.

In Figure 24.2, the mean is always 0.40 and statistically significant, but it is also important to distinguish among the three panels. In the top panel, the treatment is always moderately effective. In the middle panel, there are some populations where the treatment's effect is trivial, and some where the treatment effect is large. In the bottom panel, there are some populations where the treatment will cause minor harm, and some where it will be exceptionally helpful.

In Figure 24.3, the mean is always 0.80 and statistically significant, but it is also important to distinguish among the three panels. In the top panel, the treatment's effect is always large (as defined above). In the middle panel, there are some populations



**Figure 24.3** Three fictional examples where the mean effect is 0.80.

where the treatment's effect is moderate, and some where the treatment effect is very large. In the bottom panel, there are some populations where the treatment's effect is trivial, and some where it will be exceptionally helpful.

Unfortunately, the approach of considering the full distribution of effects is rarely taken in practice. Typically, the report of a meta-analysis will focus on the mean effect size, and then address heterogeneity as a separate matter, if at all. In fact, we need to consider the two together, and work with the entire distribution of effects, to properly address the questions we intend to address. Consider the following examples.

### METHYLPHENIDATE FOR ADULTS WITH ADHD

Castells *et al.* (2011) conducted a meta-analysis of studies that assessed the impact of methylphenidate vs. placebo on the cognitive functioning of adults with attention deficit hyperactivity disorder (ADHD). In round numbers, the mean effect size was a

standardized mean difference of roughly 0.50 with a 95% confidence interval of 0.35 to 0.65. A test of the null hypothesis (that the mean effect size is 0.0) yields a Z-value of 6.86 and a *p*-value of <0.001. The mean effect size tells us that the treatment has a moderate effect *on average*, but this is only part of the picture. To fully appreciate the potential impact of the treatment, we need to consider also the full distribution of effects, as presented in Figure 24.5.

The vertical line bounded by a confidence interval tells us that the mean effect size in the universe of comparable populations falls in the interval 0.35 to 0.65. However, in any given population, the true effect size may fall some distance from the mean. The prediction interval is roughly 0.05 to 0.95, which tells us that in some 95% of all comparable populations the true effect size will fall in this interval. The normal curve provides more detail for this distribution of effects. In most populations (within one standard deviation of the mean), the effect size is moderate. If we assume that an effect size over 0.30 is clinically important, we see that the effect is clinically important in roughly 80% of populations. Importantly, there are no populations (at least within the 95% prediction interval) where the effect is harmful. The prediction interval (labeled [B] in Figure 24.4) corresponds to the extremes of the normal curve in Figure 24.5.

Based on this information, we might suggest that the treatment can be employed immediately, but that additional research should undertaken to determine what is responsible for the variation in effect size. This suggestion is subjective, but what is clear is that any suggestion must be based on the information in Figure 24.5, and not the mean effect size alone. The suggestion would clearly be different if the dispersion in effects was more limited (on the one hand) or more expansive (on the other).

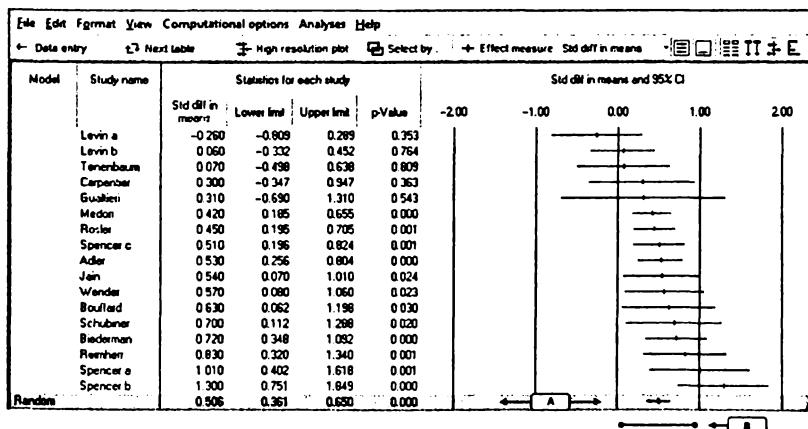
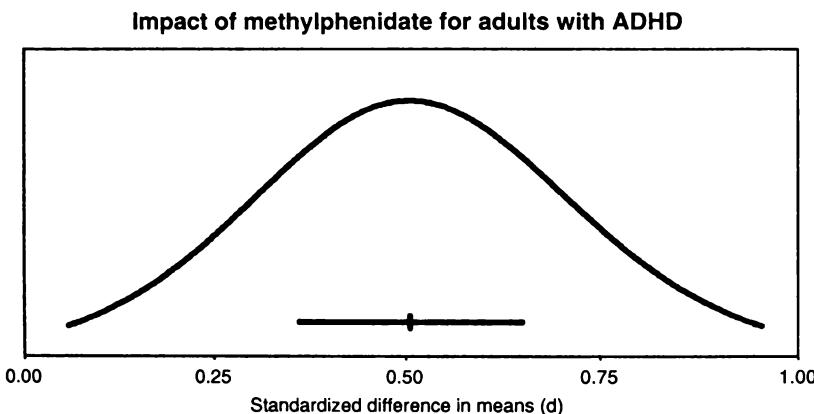


Figure 24.4 Methylphenidate for adults with ADHD (Forest plot). Effect size > 0 favors treatment.



The mean effect size is 0.51 with a 95% confidence interval of 0.36 to 0.65  
 The true effect size in 95% of all comparable populations falls in the interval 0.06 to 0.95

**Figure 24.5** Methylphenidate for adults with ADHD (True effects). Effect size > 0 favors treatment.

### IMPACT OF GLP-1 MIMETICS ON BLOOD PRESSURE

Katout *et al.* (2014) looked at the impact of GLP-1 mimetics on diastolic blood pressure (Figure 24.6). The numbers that follow are based on our reanalysis of the data and differ slightly from the original report due to rounding error.

The effect size index is the raw difference in mean blood pressure, with values below zero indicating a beneficial effect. The mean effect size is  $-0.473$ , with a confidence interval of  $-1.195$  to  $+0.248$  [B]. The confidence interval includes zero, so we *cannot* reject the null hypothesis that the mean effect size is zero. If we focused on the mean effect size, the take-home message would be that the impact of the drug, if any, is small. By contrast, if we look at the entire distribution of effects, we get a very different picture. The prediction interval [C] is roughly  $-4.08$  to  $+3.13$ . When the effects vary this widely, the mean is largely irrelevant. The prediction interval (labeled [B] in Figure 24.6) corresponds to the extremes of the normal curve in Figure 24.7.

Figure 24.7 offers a more detailed picture of the effect size distribution. The treatment may be helpful for roughly 60% of patients. However, it is harmful for the other 40%. The take-home message here should be that we need to understand where the treatment is helpful, and where it is harmful. For example, it may be helpful in specific types of patients, or in specific variants of the intervention. (The researchers did use meta-regression to try and identify potential moderators).

### AUGMENTING CLOZAPINE WITH A SECOND ANTIPSYCHOTIC

Taylor, Smith, Gee, and Nielsen (2012) looked at the impact of augmenting clozapine with a second antipsychotic (Figure 24.8). The effect size index is the standardized mean difference in function, with values below zero indicating a beneficial effect.

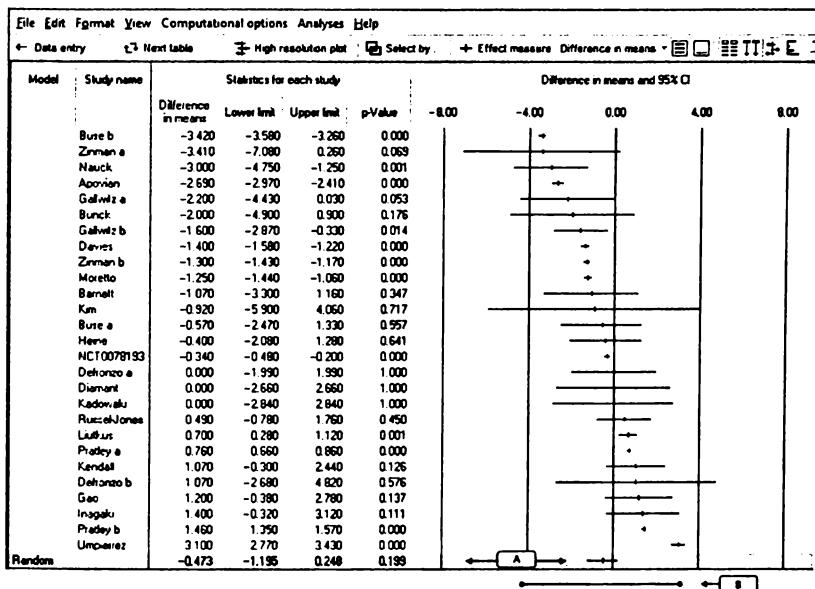


Figure 24.6 GLP-1 mimetics and diastolic BP (Forest plot). Mean difference < 0 favors treatment.

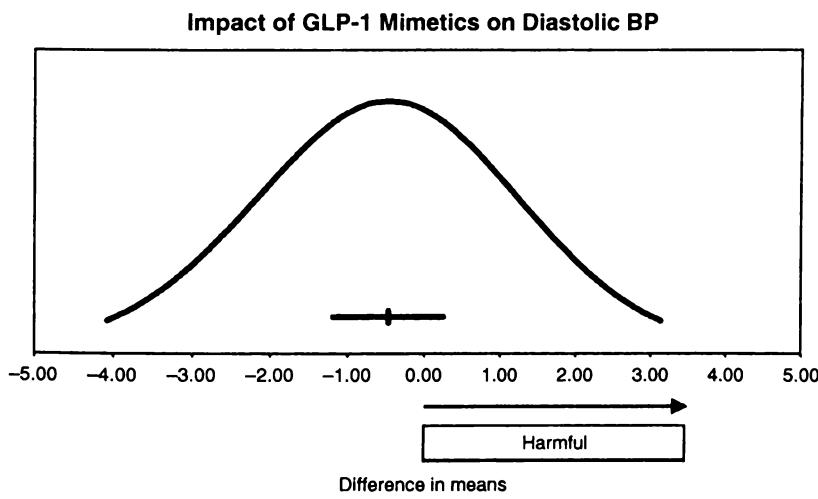


Figure 24.7 GLP-1 mimetics and diastolic BP (True effects). Mean difference < 0 favors treatment.

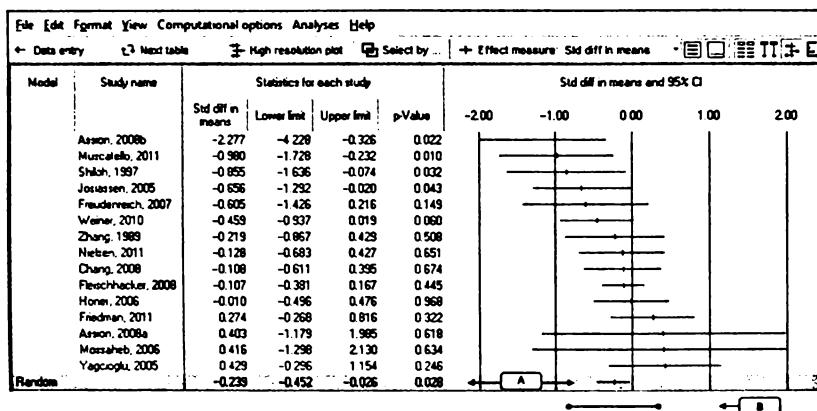
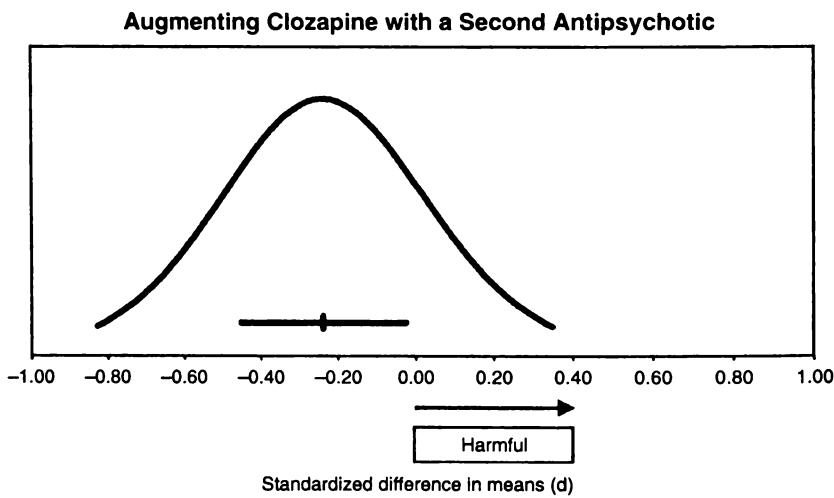


Figure 24.8 Augmenting clozapine (Forest plot). Std mean difference < 0 favors augmentation.



The mean effect size is  $-0.24$  with a 95% confidence interval of  $-0.45$  to  $-0.03$ .  
 The true effect size in 95% of all comparable populations falls in the interval  $-0.83$  to  $0.35$ .

Figure 24.9 Augmenting clozapine (True effects). Std mean difference < 0 favors augmentation.

The mean effect size is  $-0.239$  with a 95% confidence interval of  $-0.452$  to  $-0.026$  [B]. The confidence interval excludes zero, so we can reject the null hypothesis that the mean effect size is zero. If we focused on the mean effect size, the take-home message would be that the drug is helpful, on average.

By contrast, if we look at the entire distribution of effects, we get a different picture. The prediction interval (labeled [B] in Figure 24.8) corresponds to the extremes of the normal curve in Figure 24.9). The interval tells us that the effect size in any one population will be as low as -0.83 (improving function by 0.83 units) or as high as +0.35 (harming function by 0.35 standard deviations). Specifically, the treatment will be harmful in roughly 20% of populations. The take-home message here would be that we need to understand where the treatment is helpful, and where it is harmful. For example, it may be helpful in specific types of patients or in specific variants of the intervention.

## CONCLUSIONS

The key idea of a random-effects model is that the effect size will vary across studies, and we need to model this dispersion. We do model the dispersion when we compute the mean effect size and confidence interval. However, most papers do not adequately consider the implications of the dispersion when reporting the results and discussing the potential utility of the intervention.

It is not clear why this is the case. This might be a carryover from primary studies or fixed-effect meta-analyses, where we are working with one true effect size, and so there is no dispersion. It may be because the indices employed to report heterogeneity do not actually provide a clear picture of the dispersion, and so the researcher does not know what the distribution looks like. Whatever the reason, this approach is clearly not optimal. Where possible, it behooves the researcher to describe the full distribution of effects, and then consider this when discussing the potential utility of an intervention.

There is a program on the book's website that will generate the plots shown in this chapter. The program requires the user to enter the mean effect size, the upper limit of the confidence interval, the number of studies in the analysis, and  $T^2$ . The formulas used by the spreadsheet are discussed in the section on prediction intervals (Chapters 17 and 19).

## CAVEATS

It is important to keep in mind that this entire process is only useful if we have a reliable estimate of the heterogeneity. As discussed in Chapter 25 (Limitations of the random-effects analysis), this requires we meet various assumptions and that we have a sufficient number of studies in the analysis. The minimum number of studies will vary, but for a sense of scale, ten studies may serve as a useful minimum in many cases.

We should also emphasize that the distribution of effects computed here applies to the universe of comparable studies. As explained in Chapter 25, it will not always be clear what studies are included in that universe.

In the current chapter, we are focusing on the goal of understanding the dispersion in effects. This is only a first step. In the event that the effects are dispersed over a wide

area, the next logical step would be to understand why the effects vary. For example, is the treatment more effective in some populations than in others, or are some variants of the intervention more effective than others. We might be able to use subgroups analysis or regression to develop hypotheses about the impact of moderators.

### SUMMARY POINTS

- Reports of a meta-analysis typically focus on the mean effect size. The dispersion of effects is often mentioned and then ignored. While this approach is common, it is a serious mistake.
- In order to understand the potential utility of the intervention, we need to consider not only the mean effect size, but the entire distribution of effect sizes. A first step in this direction is to report the prediction interval, and then to take account of this interval when considering the potential impact of the intervention.
- A program for computing the prediction interval and plotting the distribution of effects is available on the book's website ([www.Introduction-to-Meta-Analysis.com](http://www.Introduction-to-Meta-Analysis.com)).

# Limitations of the Random-Effects Model

---

### Introduction

Assumptions of the random-effects model

A textbook case

When studies are pulled from the literature

A useful fiction

Transparency

A narrowly defined universe

Two important caveats

In context

Extreme cases

---

## INTRODUCTION

This chapter is adapted from the text *Common Mistakes in Meta-Analysis and How to Avoid Them* (Borenstein, 2019).

When the analysis is based on studies pulled from the literature, the random-effects model is almost invariably the model that should be used. This model assumes that the studies in the analysis are representative of a universe of comparable studies and that the results of the analysis will be generalized to that universe. The computation of the confidence interval and the relative weight assigned to each study reflect these goals. Critically, this model allows us to discuss not only the mean effect size, but also the dispersion in effect size across studies. These are all key goals of the analysis.

While we should be using the random-effects model for these analyses, we need to recognize that we will be violating some assumptions that are required for the model to work as intended. We need to take this into account when we interpret the results.

## ASSUMPTIONS OF THE RANDOM-EFFECTS MODEL

The random-effects model works well if the following assumptions are met.

- A. The universe to which we will be making an inference is defined clearly and is the correct universe in the sense that it is relevant to policy.
- B. The studies that were performed are a random sample from that universe.
- C. The studies that we include in our analysis are an unbiased sample of the studies that were performed.
- D. The analysis includes enough studies to yield a reliable estimate of the between-study variance,  $\tau^2$ .
- E. The true effects are normally distributed in the relevant metric.

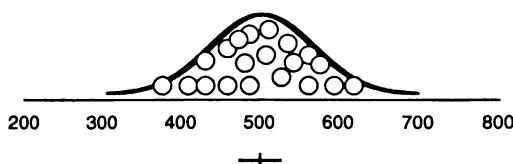
These issues build on each other. To get a reliable estimate of  $\tau^2$  in the defined universe (D), we need to have a sufficient number of cases. But we also need to assume that the studies in our analysis are a random sample of those that were performed (C), that those performed are a random sample of the defined universe (B), and that this universe is well defined and relevant to policy (A).

The quality of the evidence provided by a meta-analysis depends in large part on the extent to which that analysis meets these assumptions. If the analysis meets these assumptions fully, the quality will tend to be good. To the extent that it fails to meet some (or all) of these assumptions, the quality is likely to be poor. This list does not include assumptions about the internal validity of the studies that were performed. That is also critically important, but applies to all statistical models. For the present discussion, we assume that the individual studies have low risk of bias, and our concern is whether we can generalize from these studies to the larger universe.

## A TEXTBOOK CASE

Consider a textbook case of the random-effects model, where all assumptions of the model are fully realized. In this case, we want to estimate the mean score on a specific math test for all the high schools in New York City. We draw a random sample of 20 schools from this universe of schools and then draw a random sample of 50 students within each of these schools. This is depicted in Figure 25.1.

The 20 circles in the plot represent the true scores for the 20 schools that were included in our random sample. The key factor that makes this is a random-effects



**Figure 25.1** Random effects. Confidence interval 60 points wide.

analysis is the normal curve that has been superimposed on the plot. This curve reflects the fact that we have defined a universe of populations from which we draw the samples and to which we will be making an inference. Additionally, it is plausible to assume that the school means are normally distributed.

In this example, we can report that the statistical inference is of high quality since the assumptions have all been met. To wit:

- A. The universe to which we will be making an inference is defined as all public high schools in New York City. This is clear and unambiguous.
- B. The studies that were performed are a random sample from that universe. We know that is the case, because we had a list of all high schools in the system and used a random process to select these 20.
- C. The studies that we include in our analysis are an unbiased sample of the studies that were performed. We know that because we know that 20 studies were performed, and all 20 of them are included in our analysis.
- D. We have enough studies in our sample to yield a reliable estimate of the between-study variance.
- E. The school means are normally distributed.

## WHEN STUDIES ARE PULLED FROM THE LITERATURE

By contrast, consider what happens in a typical analysis when studies are pulled from the literature. For example, consider the analysis performed by Castells *et al.* (2011) to assess the impact of methylphenidate on adults with attention-deficit hyperactivity disorder (ADHD). Patients with this disorder have trouble performing cognitive tasks, and it was hypothesized that the drug would improve their cognitive function. The analysis includes 17 studies where patients were randomized to receive either the drug or a placebo, and then tested on measures of cognitive function. The effect size was the standardized mean difference between groups on the tests. (The original analysis includes 18 studies. Our re-analysis is based on the 17 that provided information on covariates since we use this example also for meta-regression.)

The analysis is shown in Figure 25.2, and it should be obvious that the effect size is smaller in some studies and larger in others. For purposes of this discussion, assume that the effect size tends to be lower in populations that employ a low dose of the drug and higher in populations that employ a high dose of the drug. We can use this example to highlight the differences between the textbook case and the case where studies are pulled from the literature, and show how this affects the utility of the random-effects model.

- A. We would like to think that the universe to which we are making an inference is well defined. We might think that the universe is defined adequately by the inclusion/exclusion criteria for the review, but that is rarely the case. These criteria will not entirely define the populations and methods to be included/excluded, since there are numerous factors that could influence the magnitude of the effect and we cannot enumerate all of them. Additionally, to properly define the universe we

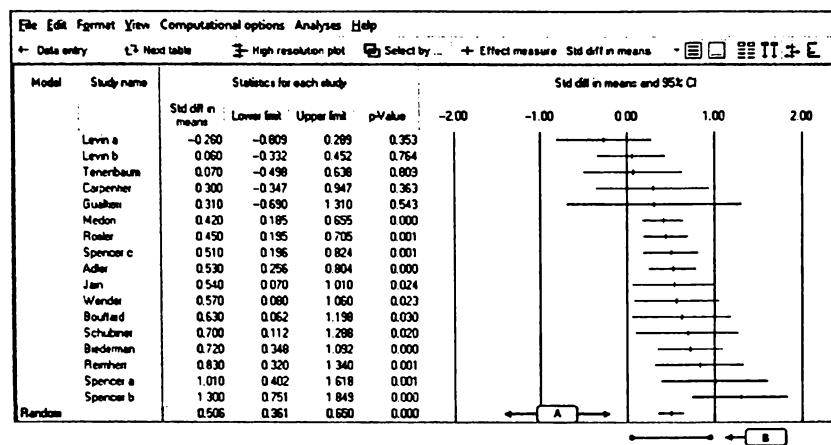


Figure 25.2 Methylphenidate for adults with ADHD. Effect size > 0 favors treatment.

would need to know not only (for example) that we will include studies where the dose is between 30 and 80 mg but also what *proportion* of studies will be using each dose in this range.

- B. We would like to think that the studies in the analysis are a *random* sample of all studies in the universe, but that is almost never the case. Researchers who perform primary studies do not design these studies using a random process. Rather, they tend to design studies that work well for their purposes and that employ populations that are relatively easy to work with. The universe defined in (A) might include equal numbers of all doses, but the studies actually performed might favor higher doses, since researchers might expect these studies tend to show larger effects.
- C. We would like to think that the studies included in the analysis are a random sample of all studies that had been performed, but for various reasons (including publication bias) the studies included in the analysis might be a biased subset of the studies that had been performed.
- D. We would like to think that we are able to estimate the between-study variance reliably, but that might not be the case. We need as many as 20 studies to obtain a reliable estimate of this variance and will often have substantially fewer studies in our analysis. Additionally, we will be estimating the between-study variance for studies in the analysis, but this may be different than the value for the studies in the *intended* universe. For example, suppose that the effect size tends to be higher for studies that employed a higher dose of the drug. If the intended universe includes all doses from 30 to 80 mg but the studies in the analysis are primarily using between 60 and 80 mg, the variance in our sample may be substantially smaller than the variance in the intended universe.

- E. When we compute the distribution of effects, we assume that we are working with a universe where the effects are normally distributed. This may not be a valid assumption.

### A USEFUL FICTION

In sum, when we apply the random-effects model to a meta-analysis where studies are pulled from the literature, we are engaging in a useful fiction. The model is useful because it provides a framework for thinking about the mean effect size and the dispersion in effects. But it is also something of a fiction because we are violating some (or all) of the assumptions that make the model work. We need to think of how these violations affect the results.

In the ADHD analysis, the mean effect size was reported as 0.51. But, given that the mean will shift left or right depending on the mix of populations in the analysis, what universe does the mean effect size represent? We cannot say that it represents the mean in the universe that we described using inclusion/exclusion rules, since we have not met assumption (A), and do not have a sampling frame. We cannot say that it represents the mean of relevant clinical populations, since we have not met assumption (B). Additionally, the mean of the observed studies might be higher than the mean of the studies actually performed since we likely violate assumption (C). In this analysis, we may have a sufficient number of studies to yield a reliable estimate of  $T^2$  for the studies *in the analysis*. Still, we may violate part of assumption (D) since the studies in the analysis may not be the same as the studies in the intended universe.

To get around these violations, we use language. We say that the results can be generalized to studies which are *comparable* to those in the analysis, without specifying what those studies are. This verbal sleight of hand yields a definition that is accurate but not useful. It is accurate since it is a tautology. The studies in the analysis are indeed comparable to studies that are comparable. But it is not useful since it does not really tell us *which* studies are comparable. That critical item is left to the judgment of the researcher or the reader, and it may not be the same as we had intended when we planned the review.

Given that there are limitations inherent in the analysis, we need to approach the results logically and see what conclusions we can draw. When we look at the entire distribution of effects, we can get a sense of the dispersion. We can then look at the mean in that context.

If there is only trivial heterogeneity among the universe of comparable studies, it follows that (a) the mean provides a useful estimate of the effect size in any given study and (b) the estimated mean will be reasonably stable regardless of which studies we happen to include in the analysis. By contrast, if there is substantial heterogeneity among the universe of comparable studies, it follows that (a) the mean does *not* provide a useful estimate of the effect size in any given study and (b) the estimated mean will vary depending on which studies we happen to include in the analysis.

For example, in the ADHD analysis there are some combinations of factors that will lead to effects as low as 0.06 and others that will lead to effects as high as 0.96.

Given the amount of heterogeneity, we should understand that the mean could shift substantially based on the particular mix of populations and methods (for example, dosage) included in the analysis. As such, the mean is not very robust.

At the same time, given the amount of heterogeneity, the mean is not terribly important. In other words, the mean is not very useful as a predictor of the effect size in any single population. Rather than focus on the mean, we need to identify factors that tell us where the effect size will be closer to 0.06 and where it will be closer to 0.96.

Since the mean itself refers to a specific (and somewhat arbitrary) mix of populations, we should recognize that the test of the null hypotheses pertains to this specific mix of populations only. For example, suppose we want to assess the impact of the drug in a universe that includes an equal number of short-term studies and long-term studies. As it happens, 80% of studies that were performed (and included in the analysis) are short-term studies. The null hypothesis we *intended* to test relates to one universe, but the null hypothesis we are *actually* testing relates to another. We might have no interest in a universe where 80% of the studies are short-term.

## TRANSPARENCY

If the violation of assumptions affects the kinds of conclusions we can draw from the analysis, we should explain what that means.

Many readers assume that the mean in the analysis pertains to the mean in some clearly designated universe. In the ADHD analysis, we should make it clear that this is not the case. The overall mean applies to the mix of populations and treatments included in the analysis, and would shift if we included a different mix of studies.

Many readers focus on the mean effect size and pay little attention to the dispersion in effects. In the ADHD analysis, we should explain that the true effect size in any given study could fall some distance from the mean. And, the mean itself could shift left or right, depending on the mix of studies included.

## A NARROWLY DEFINED UNIVERSE

In almost any meta-analysis where studies address the impact of an intervention and are pulled from the literature, we will be violating some assumptions of the random-effects model, and therefore, we need to think about the issues outlined above. However, the severity of the violations (and the potential impact) depends on several factors. Primary among these is the extent to which the universe is defined to encompass a very narrow set of studies or a more broadly defined set of studies. The ADHD analysis includes a clinically diverse set of studies and effects, but other analyses will work with a narrowly defined set of criteria. The paragraph headings below (A to D) refer to the same items listed above.

- A. When the universe is defined narrowly, it may be possible to provide a clear and comprehensive definition of the universe. Experts may be able to identify all variables that could be related to the effectiveness of the drug and set strict inclusion/exclusion criteria for these.

- B. When the universe is defined narrowly, there is less concern that the studies being performed fall toward one end or the other of a distribution, since the entire distribution is narrow. We do not need to be concerned that the dosage in our studies is higher than the typical dose in the universe since the universe is limited to one dose. The same idea applies to the type of patient, the outcome, and so on.
- C. The potential impact of publication bias depends on the extent of variation in *observed* effects, which includes variation in true effects and variation due to sampling error. When the universe is defined narrowly, bias based on variation in true effects will be limited. However, there may still be substantial bias based on sampling error.
- D. The number of studies that we need to get a reliable estimate of the between-study variance ( $\tau^2$ ) depends in part on how widely the true effects vary. When the universe is defined narrowly and the within-study variance is small, we may be able to get a reliable estimate of  $\tau^2$  with only a handful of studies.

In the Cochrane Database of Systematic Reviews, a substantial proportion of the meta-analyses report that the between-study variance is *estimated* as zero. While it is not likely that the true variance is actually *zero*, it is possible that the true variance is *trivial*. In that case, the issues outlined here for a narrowly defined universe would apply.

## TWO IMPORTANT CAVEATS

The idea that a narrowly defined universe allows us to avoid some limitations of the random-effects model depends on two important caveats.

First, the fact that  $T^2$  is estimated as zero does not tell us that the universe is narrowly defined. The estimate of  $T^2$  is typically unreliable, and even if  $T^2$  is estimated as zero, the true value could be large. Therefore, the assumption that  $T^2$  is trivial should be based primarily on logic. For example, this may be a reasonable assumption if the studies appear to be identical (or nearly so) on all factors that could be related to the effect size.

Second, if the universe is defined narrowly, then the results apply only to this narrow universe. The suggestion that the mean is reliable because we are dealing with a narrowly defined universe only applies if the results are limited to that universe. In practice, readers will often (almost invariably) extend the findings to any populations of interest. (This problem applies also to cases where the universe is defined broadly.)

## IN CONTEXT

Given the problems associated with using the random-effects model when studies are pulled from the literature, some have advocated for using the fixed-effects (plural) model instead. While there are limitations to both models, there is a growing consensus that the random-effects model is generally preferable, for the following reasons.

First, the random-effects model provides the correct conceptual framework for thinking about the analysis. It explicitly acknowledges that we intend to make an inference to a wider set of studies. Even if parts of the process are ambiguous (for example, deciding which studies are comparable), it is preferable to include them in the process so that we are clear about where the model is not reliable.

Second, the random-effects model allows us to compute prediction intervals that tell us the range of true effect sizes that might be expected in comparable studies. As explained earlier, this can be an essential element in our understanding of the results.

Third, the fixed-effects model reports a confidence interval for the mean effect size for the studies *in the analysis*, and this interval tends to be relatively narrow. By contrast, the random-effects model reports a confidence interval for the *universe of comparable studies*, and this interval tends to be wider. If we will be making an inference to the universe of comparable studies, the random-effects interval is a better match for the intended inference.

Fourth, under the fixed-effects model, large studies may dominate the analysis and small studies may be effectively ignored. Essentially, this means that we assign the same weight to each person rather than each study (see Hedges & Vevea, 1998; Peto, 1987; Rice *et al.*, 2018; “Tamoxifen for early breast cancer: an overview of the randomised trials. Early Breast Cancer Trialists’ Collaborative Group,” 1998). As such, the random-effects model is a better match if our intent is to make an inference to all comparable studies.

Fifth, while both models have limitations, the random-effects model has the *potential* to work well when we have enough data and a representative sample of the intended universe. As we approach these conditions, the model *will* yield a useful estimate of the mean effect size and dispersion of effects in that universe. By contrast, the fixed-effects model is designed to make an inference only to the studies actually included in the analysis, and not to the universe of comparable studies. This will not change as the quality and quantity of the data improve.

## EXTREME CASES

For the reasons discussed above, the random-effects model should generally be the model we use when studies are pulled from the literature. This model is likely to work *well enough* when the universe is defined narrowly. It may also work *well enough* when the universe is defined broadly but we have a reasonable number of studies, so we can get a general sense of the dispersion.

However, if the universe is defined broadly and we have only two or three studies (for example), the model becomes untenable. In this case, we need to choose among several options.

- We can apply the random-effects model and explain that the estimates are unreliable. It would be very useful to apply the Knapp–Hartung correction. This will substantially expand the width of the confidence interval and thus clarify the extent of the uncertainty. The drawback to this approach is that the interval may be so wide that we learn almost nothing from the analysis.

- We can apply the fixed-effects (plural) model and make it very clear that the results apply only to the studies in the analysis, and cannot be generalized beyond them to any other studies. The drawback to this approach is that readers will tend to ignore the caveat and generalize as they see fit. If readers do generalize, they will be working with statistics that may be dominated by the larger studies in the analysis, which generally does not make sense in this context.
- We can display the forest plot without a summary effect (Poole & Greenland, 1999). The problem with this approach is that readers may construct an even more flawed summary of their own.

None of these options is a good one, and regardless of the option chosen, it is imperative to be transparent about the limitations of the analysis.

#### SUMMARY POINTS

- When the analysis is based on studies pulled from the literature, the random-effects model is almost invariably the model that should be used. However, it is important to be aware of this model's limitations when used in this context.
- The model only yields useful estimates of the mean and heterogeneity when the analysis is based on a sufficient number of studies.
- When we do have a sufficient number of studies, we can make an inference to the universe of comparable studies, but (a) it may not be clear what that universe is and (b) that universe might not be the one that we care about.
- Since there may be no way to entirely overcome these limitations, we need to be transparent about them.



# Knapp–Hartung Adjustment

---

### Introduction

Adjustment is rarely employed in simple analyses

Adjusting the standard error

The Knapp–Hartung adjustment for other effect size indices

$t$  distribution vs.  $Z$  distribution

Limitations of the Knapp–Hartung adjustment

---

## INTRODUCTION

In meta-analysis, the confidence interval for the mean is traditionally based on the  $Z$  distribution, which yields a relatively narrow interval. When we use the random-effects model, it would be better to use the Knapp–Hartung adjustment (sometimes called the Hartung–Knapp–Sidik–Jonkman adjustment), which yields a wider (and more accurate) confidence interval (Higgins & Thompson, 2004; IntHout, Ioannidis, & Borm, 2014; Jackson, Law, Rücker, & Schwarzer, 2017; Knapp & Hartung, 2003; Sidik & Jonkman, 2002).

The adjustment includes two components. First, it modifies the standard error of the mean. Second, it multiplies the standard error by a factor based on the  $t$  distribution rather than the  $Z$  distribution. It is *always* a good idea to use this adjustment, since the adjusted interval is more accurate. However, it is *especially important* to use this adjustment when there are a small number of studies in the analysis and the between-study variance is nontrivial.

## ADJUSTMENT IS RARELY EMPLOYED IN SIMPLE ANALYSES

While there is a consensus among statisticians that we should always apply this adjustment when we use the random-effects model, it is used only rarely in practice, for several reasons.

- Most researchers are not aware of this adjustment.
- The adjustment is not always available in software.

- The adjustment may yield a very wide confidence interval and may move the  $p$ -value to a nonsignificant range, which makes it less attractive to researchers who may have a vested interest in reporting a statistically significant result.

In an effort to address the second of these items, the adjustment is being added as an option in software and hopefully will be adopted more widely in the near future. It is possible to use this adjustment now in CMA. IntHout *et al.* (2014) show how the adjustment can be implemented in Excel. These options are discussed at the end of this chapter.

## ADJUSTING THE STANDARD ERROR

IntHout *et al.* (2014) show how the adjustment can be implemented in Excel. For a motivating example we use the ADHD analysis (Castells *et al.*, 2011) which is discussed elsewhere in this volume. The data are presented in Table 26.1.

**Table 26.1** Knapp–Hartung computations for ADHD analysis.

$d$	$V$	$T^2$	$W$	$Wd$	$X$	$RES$	$RES^2 * W$
0.5300	0.0196	0.038667	17.162255	9.095995	0.505818	0.024182	0.010036
0.7200	0.0361	0.038667	13.374813	9.629865	0.505818	0.214182	0.613553
0.6300	0.0841	0.038667	8.145485	5.131655	0.505818	0.124182	0.125612
0.3000	0.1089	0.038667	6.776564	2.032969	0.505818	-0.205818	0.287063
0.3100	0.2601	0.038667	3.347085	1.037596	0.505818	-0.195818	0.128343
0.5400	0.0576	0.038667	10.387732	5.609375	0.505818	0.034182	0.012137
-0.2600	0.0784	0.038667	8.542088	-2.220943	0.505818	-0.765818	5.009745
0.0600	0.0400	0.038667	12.711746	0.762705	0.505818	-0.445818	2.526511
0.4200	0.0144	0.038667	18.843960	7.914463	0.505818	-0.085818	0.138782
0.8300	0.0676	0.038667	9.410223	7.810485	0.505818	0.324182	0.988955
0.4500	0.0169	0.038667	17.996163	8.098273	0.505818	-0.055818	0.056070
0.7000	0.0900	0.038667	7.771976	5.440383	0.505818	0.194182	0.293054
1.0100	0.0961	0.038667	7.420192	7.494394	0.505818	0.504182	1.886206
1.3000	0.0784	0.038667	8.542088	11.104714	0.505818	0.794182	5.387704
0.5100	0.0256	0.038667	15.559988	7.935594	0.505818	0.004182	0.000272
0.0700	0.0841	0.038667	8.145485	0.570184	0.505818	-0.435818	1.547134
0.5700	0.0625	0.038667	9.884607	5.634226	0.505818	0.064182	0.040717
SUM			184.022449	93.081935			19.051895

where

- $d$  Standardized mean difference  
 $V_d$  Within-study error variance  
 $T^2$  Tau-squared (between-study variance)  
 $W$  Study weight  $1/(V+T^2)$   
 $Wd$  Product of  $d$  times the weight  
 $X$  Mean effect size using random-effects weights  
 $Res$  Residual,  $d-X$   
 $Res^2 * W$  Residual squared times  $W$

Without the adjustment, the results would be computed as follows.

$$M = \frac{\sum Wd}{\sum W} = \frac{93.081935}{184.022449} = 0.505818 \quad (26.1)$$

$$V_M = \frac{1}{\sum W_i} = \frac{1}{184.022449} = 0.005434 \quad (26.2)$$

$$SE_M = \sqrt{V_M} = \sqrt{0.005434} = 0.073716 \quad (26.3)$$

$$Z_{CRIT} = ABS(NORM.S.INV(0.05/2)) = 1.959964 \quad (26.4)$$

$$CI_{LL} = M - Z_{CRIT} \times SE = 0.505818 - 1.959964 \times 0.073716 = 0.361337 \quad (26.5)$$

$$CI_{UL} = M + Z_{CRIT} \times SE = 0.505818 + 1.959964 \times 0.073716 = 0.650300 \quad (26.6)$$

To test the null hypothesis that  $M = 0.00$  we compute

$$Z = \frac{M}{SE_M} = \frac{0.505818}{0.073716} = 6.861673 \quad (26.7)$$

$$p = (1-NORMSDIST(6.861673)) * 2 = 0.000000 \quad (26.8)$$

where

$M$  Mean effect

$V_M$  Variance of mean

$SE_M$  Standard error of mean

$Z_{CRIT}$  Criterion Z value

$CI_{LL}$  Confidence interval lower limit

$CI_{UL}$  Confidence interval upper limit

$Z$  Z value for test of null hypothesis

$p$  p-value for test of null hypothesis

With the adjustment, the results would be computed as follows.

$$M = \frac{\sum Wd}{\sum W} = \frac{93.081935}{184.022449} = 0.505818 \quad (26.9)$$

$$Factor = \frac{\sum Res^2 \times W}{\sum W} = \frac{19.051895}{184.022449} = 0.103530 \quad (26.10)$$

$$V_M = \frac{Factor}{df} = \frac{0.103530}{16} = 0.006471 \quad (26.11)$$

$$SE_M = \sqrt{V_M} = \sqrt{0.006471} = 0.080440 \quad (26.12)$$

$$t_{(CRIT)} = TINV(0.05, 16) = 2.119905 \quad (26.13)$$

$$CI_{LL} = M - t(SE) = 0.505818 - 2.119905 \times 0.080440 = 0.335293 \quad (26.14)$$

$$CI_{UL} = M + t(SE) = 0.505818 + 2.119905 \times 0.080440 = 0.676344 \quad (26.15)$$

To test the null hypothesis that  $M = 0.00$ , we compute

$$t = \frac{M}{SE_M} = \frac{0.505818}{0.080440} = 6.288121 \quad (26.16)$$

$$p = TDIST(T, df, tails) = TDIST(6.288121, 16, 2) = 0.000011 \quad (26.17)$$

where

$M$	Mean effect
$V_M$	Variance of mean
$SE_M$	Standard error of mean
$t_{CRT}$	Criterion $t$ value
$CI_{LL}$	Confidence interval lower limit
$CI_{UL}$	Confidence interval upper limit
$t$	$t$ value for test of null hypothesis
$p$	$p$ -value for test of null hypothesis

Table 26.2 shows the adjusted values alongside side the original values. In this example, the adjusted values are fairly close to the original values, but this is because the analysis includes seventeen studies and the estimate of the mean is reasonably precise. In other cases, the adjustment will have a larger impact. In analyses with only a few studies, the adjustment can increase the width of the confidence interval by twofold more.

Table 26.2 Original vs. Knapp–Hartung.

	Original	Knapp–Hartung
$M$	0.505818	0.505818
$V_M$	0.005434	0.006471
$SE_M$	0.073716	0.080440
Criterion	1.959964	2.119905
$CI_{LL}$	0.361337	0.335293
$CI_{UL}$	0.650300	0.676344
Test statistic	6.861673	6.288121
$p$ -value	0.000000	0.000011

The formulas in (26.4), (26.8), (26.13), and (26.17) are for the version of Excel™ distributed in 2020.

### THE KNAPP–HARTUNG ADJUSTMENT FOR OTHER EFFECT SIZE INDICES

When the effect size index is  $d$ ,  $g$ ,  $D$ , or  $RD$ , the analyses are performed in the same metric as the effect size itself and so (as in this example) the adjustments are based on the raw data. When the effect size is a ratio, the analyses are performed using log

values, and the adjustments would be performed using log values. In that case, the  $t$  value and corresponding  $p$ -value are valid as computed using (26.16) and (26.17). However, the confidence interval in (26.14) and (26.15) is computed in log units and must be converted to ratio units. For the risk ratio, we would use

$$RR = \exp(\log_{RR}). \quad (26.18)$$

When the effect size is a correlation, the analyses may be performed using Fisher's  $Z$  values. In that case, the adjustments would also be performed using Fisher's  $Z$  values. The  $t$  value and corresponding  $p$ -value are valid as computed using (26.16) and (26.17). However, the confidence interval in (26.14) and (26.15) is computed in Fisher's  $Z$  units and must be converted to correlations using

$$r = \frac{\exp(2^*Z) - 1}{\exp(2^*Z) + 1} \quad (26.19)$$

When the effect size is risk in one group, the analyses may be performed using the logit transformation. In that case, the adjustments would also be performed using the logit transformation. The  $t$  value and corresponding  $p$ -value are valid as computed using (26.16) and (26.17). However, the confidence interval in (26.14) and (26.15) is computed in logit units and must be converted to risk using

$$\text{Risk} = \frac{\exp(\text{Logit})}{\exp(\text{Logit}) + 1} \quad (26.20)$$

CMA software will include the Knapp–Hartung adjustment as part of all analyses in the near future. In the third version, it includes this adjustment only in the regression module. To apply the adjustment to a simple analysis, run that analysis as a regression with only the intercept and choose the Knapp–Hartung option. The ADHD analysis yields the following results, which match the results displayed above.

---

Random-effects, Knapp–Hartung, standardized difference in means						
Covariate	Coefficient	Standard error	95% lower	95% upper	$t$ value $df = 16$	2-sided $p$ -value
Intercept	0.505818	0.080440	0.335293	0.676344	6.288121	0.000011

---

If this procedure is employed for ratios, correlations, or prevalence, the  $t$  value and  $p$ -value are valid as displayed, but the confidence limits must be converted using (26.18), (26.19), or (26.20).

While the computations shown above are based on the formulas in IntHout *et al.* (2014), the tables in that paper include some errors. Therefore, the results shown here will not match those which follow the example in that paper.

## **t DISTRIBUTION VS. Z DISTRIBUTION**

Consider an analysis where the effect size index is the standardized mean difference ( $d$ ). The mean effect is 0.50, and the standard error of the mean is 0.10. Table 26.3

**Table 26.3** Impact of using *t* distribution on the confidence interval width.

	Number of studies	Critical value	Lower limit	Upper limit	Width	Ratio <i>t</i> : <i>Z</i>
Z distribution	<i>n/a</i>	1.960	0.304	0.696	0.392	1.00
<i>t</i> distribution	100	1.984	0.302	0.698	0.397	1.01
	30	2.045	0.295	0.705	0.409	1.04
	20	2.093	0.291	0.709	0.419	1.07
	10	2.262	0.274	0.726	0.452	1.15
	5	2.776	0.222	0.778	0.555	1.42
	4	3.182	0.182	0.818	0.636	1.62
	3	4.303	0.070	0.930	0.861	2.20
	2	12.706	-0.771	1.771	2.541	6.48

shows how the confidence interval is affected when we use *t* rather than *Z*. Without the correction, the confidence interval width is around 0.40 regardless of the number of studies. With the correction, the width increases as the number of studies decreases. When the number of studies is 30, 10, 4, and 2, the interval width is approximately 0.41, 0.45, 0.64, and 2.54. Equivalently, the width is increased by a factor of 1.04, 1.15, 1.62, and 6.48. As noted above, there is a second part to the adjustment which involves the standard error, and which may widen the interval even further.

While this discussion has been focused on the width of the confidence interval, the same issues apply to tests of the null hypothesis.

## LIMITATIONS OF THE KNAPP–HARTUNG ADJUSTMENT

For the random-effects model, the Knapp–Hartung adjustment always yields better coverage than the nonadjusted value, and so it should always be used. However, it works better under some circumstances than others (IntHout *et al.*, 2014).

The adjustment makes it more likely that the confidence interval will include the true mean for studies comparable to those in the analysis. However, it cannot adjust for the possibility that the studies in the analysis are not representative of the intended universe. For example, suppose that the universe is defined as studies that employed a dose between 30 and 80 mg, and the true mean effect size for these studies is 0.50. However, most studies in the analysis employed a dose between 60 and 80 mg, and the true mean effect size for these studies is 0.70. The Knapp–Hartung adjustment makes it more likely that the confidence interval will include the value of 0.70. It cannot adjust for the fact that this is different than the mean in the *intended* universe.

Ironically, when the adjustment is most needed (when we have a small number of studies) the impact of the adjustment may be so large that the estimate of the mean effect size will be uninformative. In the example presented in Table 26.3, with 20 studies the interval would have a width of 0.42 but with two studies would have a width of 2.54. The latter interval is so wide that it tells us nothing of real value. While this is unfortunate, it represents the true state of affairs. When the between-study variance is nontrivial, an estimate of the mean effect size for the universe of comparable studies, based on two studies, is not reliable.

**SUMMARY POINTS**

- The confidence interval for the mean effect size in random-effects analysis tends to be too narrow when based on the Z distribution. It would be better to use the Knapp–Hartung adjustment, which yields a wider interval. The adjustment applies both to the confidence interval and to the test of the null hypothesis for the mean effect.
- The magnitude of the adjustment depends on the number of studies in the analysis. When the analysis includes many studies, the adjustment will tend to be relatively modest. When the analysis includes only a few studies, the adjustment will tend to be substantial.



# Complex Data Structures



# Overview

Thus far we have assumed that each study contributes one (and only one) effect size to a meta-analysis. In this section we consider cases where studies contribute more than one effect size to the meta-analysis. These usually fall into one of the following types.

### Multiple independent subgroups within a study

Sometimes a single study will report data for several cohorts of participants. For example, if researchers anticipate that the treatment effect (e.g., drug versus placebo) could vary by age, they might report the treatment effect separately for children and for adults. Similarly, if researchers anticipate that the treatment effect could vary by disease stage they might report the effect separately for patients enrolled with early-stage disease and for those enrolled with late-stage disease.

The defining feature here is that the subgroups are independent of each other, so that each provides unique information. For this reason, we can treat each subgroup as though it were a separate study, which is sometimes the preferred method. However, there are sometimes other options to consider, and we will be discussing these as well.

### Multiple outcomes or time-points within a study

In some cases researchers will report data on several related, but distinct outcomes. A study that looked at the impact of tutoring might report data on math scores and also on reading scores. A study that looked at the association between diet and cardiovascular disease might report data on stroke and also on myocardial infarction. Similarly, a study that followed subjects over a period of time may report data using the same scale but at a series of distinct time-points. For example, studies that looked at the impact of an intervention to address a phobia might collect data at one month, six months, and twelve months.

The defining feature here is that *the same* participants provide data for the different outcomes (or time-points). We cannot treat the different outcomes as though they were independent as this would lead to incorrect estimates of the variance for the summary effect. We will show how to correct the variance to take account of the relationship among the outcomes.

## More than one comparison group within a study

Sometimes, a study will include several treatment groups and a single control group. For example, one effect size may be defined as the difference between the placebo group and drug *A*, while another is defined as the difference between the same placebo group and drug *B*.

The defining feature here is similar to multiple outcomes, in that some participants (those in the control group) contribute information to more than one effect size. The methods proposed for dealing with this problem are similar to those proposed for multiple outcomes. They also include some options that are unique to the case of multiple comparisons.

## How this Part is organized

The next three chapters address each of these cases in sequence. Within each chapter we first show how to combine data to yield a *summary* effect, and then show how to look at *differences* in effects.

The worked examples in these chapters use the fixed-effect model. We adopt this approach because it involves fewer steps and thus allows us to focus on the issue at hand, which is how to compute an effect size and variance. Once we have these effect sizes we can use them for a fixed-effect or a random-effects analysis, and the latter is generally more appropriate.

In the worked examples we deliberately use a generic effect size rather than specifying a particular effect size such as a standardized mean difference or a log odds ratio. The methods discussed here can be applied to any effect size, including those based on continuous, binary, correlational, or other kinds of data. As always, computations for risk ratios or odds ratios would be performed using log values, and computations for correlations would be performed using Fisher's z transformed values.

# Independent Subgroups within a Study

---

### Introduction

Combining across subgroups

Comparing subgroups

---

## INTRODUCTION

The first case of a complex data structure is the case where studies report data from two or more independent subgroups.

Suppose we have five studies that assessed the impact of a treatment on a specific type of cancer. All studies followed the same design, with patients randomly assigned to either standard or aggressive treatment for two months. In each study, the results were reported separately for patients enrolled with stage-1 cancer and for those enrolled with stage-2 cancer. The stage-1 and stage-2 patients represent two independent subgroups since each patient is included in one group or the other, but not both.

If our goal was to compute the summary treatment effect for all stage-1 patients and, separately, for all stage-2 patients, then we would perform two separate analyses. In this case we would treat each subgroup as a separate study, and include the stage-1 studies in one analysis and the stage-2 studies in the other.

This chapter addresses the case where we want to use data from two or more subgroups in the *same* analysis. Specifically,

- We want to compute a summary effect for the impact of the intervention for stage-1 and stage-2 patients combined.
- Or, we want to compare the effect size for stage-1 patients versus stage-2 patients.

## COMBINING ACROSS SUBGROUPS

The defining feature of independent subgroups is that each subgroup contributes independent information to the analysis. If the sample size within each subgroup is 100, then the effective sample size across two subgroups is 200, and this will be reflected in

**Table 28.1** Independent subgroups – five fictional studies.

Study	Subgroup	ES	Variance
Study 1	Stage 1	0.300	0.050
	Stage 2	0.100	0.050
Study 2	Stage 1	0.200	0.020
	Stage 2	0.100	0.020
Study 3	Stage 1	0.400	0.050
	Stage 2	0.200	0.050
Study 4	Stage 1	0.200	0.010
	Stage 2	0.100	0.010
Study 5	Stage 1	0.400	0.060
	Stage 2	0.300	0.060

the precision of the summary effect. However, within this framework we have several options for computing the summary effect.

We shall pursue the example of five studies that report data separately for patients enrolled with stage-1 or stage-2 cancer. The effect size and variance for each subgroup are shown in Table 28.1 and are labeled simply *ES* and *Variance* to emphasize the point that these procedures can be used with any effect size. If the outcome was continuous (means and standard deviations), the effect size might be a standardized mean difference. If the outcome was binary (for example, whether or not the cancer had metastasized), the effect size might be the log risk ratio.

### Using subgroup as unit of analysis (option 1a)

One option is simply to treat each subgroup as a separate study. This is shown in Table 28.2, where each subgroup appears on its own row and values are summed across the ten rows.

Then, using formulas (11.3) to (11.5),

$$M = \frac{76.666}{413.333} = 0.1855,$$

with variance

$$V_M = \frac{1}{413.333} = 0.0024$$

and standard error

$$SE_M = \sqrt{0.0024} = 0.0492.$$

### Using study as unit of analysis (option 1b)

A second option is to compute a composite score for each study and use this in the analysis, as in Figure 28.1. The unit of analysis is then the study rather than the subgroup.

**Table 28.2** Independent subgroups – summary effect.

Study	Subgroup	ES	Variance	WT	ES*WT
Study 1	Stage 1	0.30	0.05	20.000	6.000
	Stage 2	0.10	0.05	20.000	2.000
Study 2	Stage 1	0.20	0.02	50.000	10.000
	Stage 2	0.10	0.02	50.000	5.000
Study 3	Stage 1	0.40	0.05	20.000	8.000
	Stage 2	0.20	0.05	20.000	4.000
Study 4	Stage 1	0.20	0.01	100.000	20.000
	Stage 2	0.10	0.01	100.000	10.000
Study 5	Stage 1	0.40	0.06	16.667	6.667
	Stage 2	0.30	0.06	16.667	5.000
Sum				413.333	76.667

**Computing a combined effect across subgroups within a study**

The mean and variance of the composite within a study are computed by performing a fixed-effect meta-analysis on the subgroups for that study. For study 1, this is shown in Figure 28.1 and in Table 28.3.

We apply formulas (11.3) and (11.4) to yield a mean effect

$$M = \frac{8.0000}{40.0000} = 0.2000$$

with variance

$$V_M = \frac{1}{40.0000} = 0.0250.$$

Study	Subgroup	Effect size	Variance	Mean	Variance
Study 1	Stage-1	0.30	0.05	0.200	0.025
	Stage-2	0.10	0.05		

**Figure 28.1** Creating a synthetic variable from independent subgroups.**Table 28.3** Independent subgroups – synthetic effect for study 1.

Subgroup	Effect Y	Variance V <sub>Y</sub>	Weight W	Computed WY
Stage 1	0.30	0.05	20.000	6.000
Stage 2	0.10	0.05	20.000	2.000
Sum			40.000	8.000

**Table 28.4** Independent subgroups – summary effect across studies.

Study	Subgroup	ES	Variance	ES	Variance	WT	ES*WT
Study 1	Stage 1	0.300	0.050	0.200	0.025	40.000	8.000
	Stage 2	0.100	0.050				
Study 2	Stage 1	0.200	0.020	0.150	0.010	100.000	15.000
	Stage 2	0.100	0.020				
Study 3	Stage 1	0.400	0.050	0.300	0.025	40.000	12.000
	Stage 2	0.200	0.050				
Study 4	Stage 1	0.200	0.010	0.150	0.005	200.000	30.000
	Stage 2	0.100	0.010				
Study 5	Stage 1	0.400	0.060	0.350	0.030	33.333	11.667
	Stage 2	0.300	0.060				
Sum						413.333	76.667

Note that the variance for the study (0.025) is one-half as large as the variance for either subgroup (0.050) since it is based on twice as much information.

This procedure is used to form a composite effect size and variance for each study, as shown in Table 28.4. Then, we perform a meta-analysis working solely with these study-level effect sizes and variances.

At this point we can proceed to the meta-analysis using these five (synthetic) scores. To compute a summary effect and other statistics using the fixed-effect model, we apply the formulas starting with (11.3). Using values from the line labeled *Sum* in Table 28.4,

$$M = \frac{76.667}{413.333} = 0.1855$$

with variance

$$V_M = \frac{1}{413.333} = 0.0024$$

and standard error

$$SE_M = \sqrt{0.0024} = 0.0492.$$

Note that the summary effect and variance computed using study as the unit of analysis are identical to those computed using subgroup as the unit of analysis. This will always be the case when we use a fixed effect analysis to combine effects at both steps in the analysis (within studies and across studies).

However, the two methods will yield different results if we use a random effects analysis to combine effects across studies. This follows from the fact that  $T^2$  (which is used to compute weights) may be different if based on variation in effects from study to study than if based on variation in effects from subgroup to subgroup. Therefore, the decision to use subgroup or study as the unit of analysis should be based on the context for computing  $T^2$ . Consider the following two cases.

- Case 1: Five researchers studied the impact of an intervention. Each researcher selected five schools at random, and each school is reported as an independent subgroup within the study. In this situation, between-school variation applies just as much within studies as across studies. We would therefore use subgroup as the unit of analysis (option 1a).

- Case 2: Five researchers have published papers on the impact of an intervention. Each researcher worked in a single school, and the subgroups are grade 1, grade 2, and so on. We expect the effects to be relatively consistent within a school, but to vary substantially from one school to the next. To allow for this between-school variation we should use a random-effects model only *across studies*. To properly estimate this component of uncertainty we would use study as the unit of analysis (option 1b).

### Recreating the summary data for the full study (option 2)

Options 1a and 1b differ in the unit of analysis, but they have in common that the effect size is computed *within* subgroups. Another option (option 2) is to use the summary data from the subgroups to recreate the data for the study as a whole, and then use this summary data to compute the effect size and variance.

When the subgroup data are reported as  $2 \times 2$  tables we can simply collapse cells to recreate the data for the full sample. That is, we sum cell A over all subgroups to yield an overall cell A, and repeat the process for cells B, C, and D.

When the subgroup data are reported as means, standard deviations, and sample size for each treatment group the combined sample size is summed across subgroup. For example, for treatment group 1,

$$n_1 = n_{11} + n_{12}, \quad (28.1)$$

the combined mean is computed as the weighted mean (by sample size) across groups,

$$\bar{X}_1 = \frac{n_{11}\bar{X}_{11} + n_{12}\bar{X}_{12}}{n_{11} + n_{12}}, \quad (28.2)$$

and the combined standard deviation is computed as

$$S_1 = \sqrt{\frac{(n_{11} - 1)S_{11}^2 + (n_{12} - 1)S_{12}^2 + \frac{n_{11}n_{12}}{n_{11} + n_{12}}(\bar{X}_{11} - \bar{X}_{12})^2}{n_{11} + n_{12} - 1}}, \quad (28.3)$$

where  $\bar{X}_{11}$ ,  $\bar{X}_{12}$  are the means in subgroups 1 and 2 of treatment group 1;  $S_{11}$ ,  $S_{12}$  the standard deviations, and  $n_{11}$ ,  $n_{12}$  the sample sizes; of subgroups 1 and 2.

When the subgroup data are reported as correlations, analogous formulas exist to recreate the correlation for the full study, but these are beyond the scope of this book.

Option 2 is sometimes used when some studies report summary data for all subjects combined, while others break down the data by subgroups. If the researcher believes that the subgroup classifications are unimportant, and wants to have a uniform approach for all studies (to compute an effect size from a single set of summary data) then this option will prove useful.

However, it is important to understand that this is a fundamentally different approach than other options. To return to the example introduced at the start of this chapter, under options 1a and 1b the effect size was computed *within* subgroups, which means that the effect size is *the impact of intervention controlling for cancer stage* (even if we then merge the effect sizes to yield an overall effect). By contrast, under option 2 we merge the summary data and *then* compute an effect size. Therefore the effect is ‘*the impact of intervention ignoring cancer stage*’.

When the studies are randomized trials, the proportion of participants assigned to each treatment is typically constant from one subgroup to the next. In this case there is not likely to be a confounder between treatment and subgroup, and so either approach would be valid. By contrast, in observational studies the proportion of exposed subjects may vary from one subgroup to the next, which would yield confounding between exposure and subgroup. In this case option 2 should *not* be used. (This is the same issue discussed in Chapter 38, under the heading of Simpson's paradox.)

## COMPARING SUBGROUPS

When our goal is to *compare* the effect size in different subgroups (rather than *combine* them) we have two options, as follows.

### Using subgroup as unit of analysis

One option is simply to treat each subgroup as a separate study, where each study is classified (in this example) as stage 1 or stage 2. We then compute a summary effect for all the stage 1 effects, another for all the stage 2 effects, and then compare the two using a Z-test or analysis of variance as discussed in Chapter 21.

A second option is to compute the effect size within subgroups for each study, and then to compute the difference in effects within each study. In this case each study will contribute one effect to the analysis, where the effect is the difference between subgroups.

The first option is a more general approach. It allows us to work with studies that report data for any subgroup or combination of subgroups (one study has subgroups A and B, another B and C, and so on), and then to use all relevant subgroups to compute the summary effect.

The second option can only be used if all studies report data on the same two subgroups, which is relatively rare. When this option *can* be used, however, it will usually yield a more precise estimate of the difference in effects in random effects analyses, and is also desirable because differences in effects are not confounded by possible differences between studies.

### SUMMARY POINTS

- When we have independent subgroups within a study, each subgroup contributes independent information. Therefore, if the sample size within each subgroup is 100, then the effective sample size across five subgroups is 500. In this sense, independent subgroups are no different than independent studies.
- To compute a summary effect we typically compute the effect within subgroups and then either use these effects as the unit of analysis, or merge effects within

each study and use study as the unit of analysis. A second option is to combine the summary data from all subgroups to recreate the original study level data, and then compute an effect size from this data. The second approach should be used only in limited circumstances.

- To compare effects across subgroups we typically use subgroup as the unit of analysis. In some cases we may also be able to compute the difference between subgroups in each study, and use study as the unit of analysis.



# Multiple Outcomes or Time-Points within a Study

---

### Introduction

Combining across outcomes or time-points

Comparing outcomes or time-points within a study

---

## INTRODUCTION

The second case of a complex data structure is the case where a study reports data on more than one outcome, or more than one time-point, where the different outcomes (or time-points) are based on the same participants.

For example, suppose that five studies assessed the impact of tutoring on student performance. All studies followed the same design, with students randomly assigned to either of two groups (tutoring or control) for a semester, after which they were tested for proficiency in reading and math. The effect was reported separately for the reading and the math scores, but within each study *both outcomes were based on the same students*.

Or, consider the same situation with the following difference. This time, assume that each study tests only for reading but does so at two time-points (immediately after the intervention and again six months later). The effect was reported separately for each time-point but *both measures were based on the same students*.

For our purposes the two situations (multiple outcomes for the same subjects or multiple time-points for the same subjects) are identical, and we shall treat them as such in this discussion. We shall use the term *outcomes* throughout this chapter, but the reader can substitute *time-points* in every instance.

If our goal was to compute a summary effect for the impact of the intervention on reading, and *separately* for the impact of the intervention on math scores, we would simply perform two separate meta-analyses, one using the data for reading and the other using the data for math. The issues we address in this chapter are how to proceed when we want to incorporate both outcomes in the same analysis. Specifically,

- We want to compute a summary effect for the intervention on *Basic skills*, which combines the data from reading and math.
- Or, we want to investigate the *difference* in effect size for reading versus math.

In either case, the issue we need to address is that the data for reading and math are not independent of each other and therefore the errors are correlated.

## COMBINING ACROSS OUTCOMES OR TIME-POINTS

The data for the five fictional studies are shown in Table 29.1. In study 1, for example, the effect size for reading was 0.30 with a variance of 0.05, and the effect size for math was 0.10 with a variance of 0.05.

While it might seem that we could treat each line of data as a separate study and perform a meta-analysis with ten *studies*, this is problematic for two reasons. One problem is that in computing the *summary* effect across studies this approach will assign more weight to studies with two outcomes than to studies with one outcome. (While this problem does not exist in our set of studies, it would be a problem if the number of outcomes varied from study to study.)

The second, and more fundamental problem, is that this approach leads to an improper estimate of the precision of the summary effect. This is because it treats the separate outcomes as providing independent information, when in fact the math and reading scores come from the same set of students and therefore are not independent of each other. If the outcomes are positively correlated (which is almost always the case with effects that we would want to combine), this approach underestimates the error (and overestimates the precision) of the summary effect.

Note. If the correlation between outcomes is negative, this approach will overestimate the error (and underestimate the precision) of the summary effect. The solutions presented below will work for this case as well, but in the discussion we assume that we are dealing with a positive correlation.

To address these problems, rather than treating each outcome as a separate unit in the analysis, we'll compute the mean of the outcomes for each study, and use this synthetic score as the unit of analysis. In Table 29.2 we show this schematically for study 1.

**Table 29.1** Multiple outcomes – five fictional studies.

Study	Outcome	ES	Variance
Study 1	Reading	0.300	0.050
	Math	0.100	0.050
Study 2	Reading	0.200	0.020
	Math	0.100	0.020
Study 3	Reading	0.400	0.050
	Math	0.200	0.050
Study 4	Reading	0.200	0.010
	Math	0.100	0.010
Study 5	Reading	0.400	0.060
	Math	0.300	0.060

**Table 29.2** Creating a synthetic variable as the mean of two outcomes.

Study	Outcome	Effect size	Variance	Mean	Variance
Study 1	Math	0.30	0.05	0.20	?
	Reading	0.10	0.05		

We start with summary data for two outcomes (math and reading), and compute an effect size and variance for each. If the data are continuous (means and standard deviations on the exam) the effect size might be Hedges'  $g$ . If the data are binary (number of students passing the course) the effect size might be a log risk ratio. And so on. Then, we compute a synthetic effect size for *Basic skills* which incorporates both the math and reading effects. The method used to compute this effect size and its variance is explained below.

Since every study will be represented by one score in the meta-analysis regardless of the number of outcomes included in the mean, this approach solves the problem of more weight being assigned to studies with more outcomes. This approach also allows us to address the problem of non-independent information, since the formula for the variance of the synthetic variable will take into account the correlation among the outcomes.

### **Computing a combined effect across outcomes**

Our notation will be to use  $Y_1$ ,  $Y_2$  etc. for effect sizes from different outcomes or time points within a study, and  $Y_j$  to refer to the  $j^{\text{th}}$  of these. Strictly, we should use  $Y_{ij}$ , for the  $j^{\text{th}}$  outcome (or time-point) in the  $i^{\text{th}}$  study. However, we drop the  $i$  subscript for convenience. The effect size for *Basic skills* is computed as the mean of the reading and math scores,

$$\bar{Y} = \frac{1}{2}(Y_1 + Y_2). \quad (29.1)$$

This is what we would use as the effect estimate from this study in a meta-analysis. Using formulas described in Box 29.1, the variance of this mean is

$$V_{\bar{Y}} = \frac{1}{4}(V_{Y_1} + V_{Y_2} + 2r\sqrt{V_{Y_1}}\sqrt{V_{Y_2}}) \quad (29.2)$$

### **BOX 29.1 COMPUTING THE VARIANCE OF A COMPOSITE OR A DIFFERENCE**

#### *1. The variance of the sum of two correlated variables*

If we know that the variance of  $Y_1$  is  $V_1$  and the variance of  $Y_2$  is  $V_2$ , then

$$\text{var}(Y_1 + Y_2) = V_1 + V_2 + 2r\sqrt{V_1}\sqrt{V_2},$$

where  $r$  is the correlation coefficient that describes the extent to which  $Y_1$  and  $Y_2$  co-vary. If  $Y_1$  and  $Y_2$  are inextricably linked (so that a change in one determines completely the change in the other), then  $r = 1$ , and the variance of the sum is roughly twice the sum of the variances. At the other extreme, if  $Y_1$  and  $Y_2$  are

**BOX 29.1 CONTINUED**

unrelated, then  $r = 0$  and the variance is just the sum of the individual variances. This is because when the variables are unrelated, knowing both gives us twice as much information, and so the variance is halved compared with the earlier case.

**2. The impact of a scaling factor on the variance**

If we know the variance of  $X$ , then the variance of a scalar (say  $c$ ) multiplied by  $X$  is given by

$$\text{var}(cX) = c^2 \times \text{var}(X).$$

**3. The variance of the mean of two correlated variables**

Combining 1 with 2, we can see that the variance of the mean of  $Y_1$  and  $Y_2$  is

$$\text{var}\left(\frac{1}{2}(Y_1 + Y_2)\right) = \left(\frac{1}{2}\right)^2 \text{var}(Y_1 + Y_2) = \frac{1}{4}(V_1 + V_2 + 2r\sqrt{V_1}\sqrt{V_2}).$$

**4. The variance of the sum of several correlated variables**

If we know  $Y_i$  has variance  $V_i$  for several variables  $i = 1, \dots, m$ , then the formula in 1 extends as follows:

$$\text{var}\left(\sum_{i=1}^m Y_i\right) = \sum_{i=1}^m V_i + \sum_{i \neq j} (r_{ij}\sqrt{V_i}\sqrt{V_j})$$

where  $r_{ij}$  is the correlation between  $Y_i$  and  $Y_j$ .

**5. The variance of the mean of several correlated variables**

Combining 4 with 2, we can see that the variance of the mean of several variables is

$$\text{var}\left(\frac{1}{m} \sum_{i=1}^m Y_i\right) = \left(\frac{1}{m}\right)^2 \text{var}\left(\sum_{i=1}^m Y_i\right) = \left(\frac{1}{m}\right)^2 \left(\sum_{i=1}^m V_i + \sum_{i \neq j} (r_{ij}\sqrt{V_i}\sqrt{V_j})\right).$$

**6. The variance of the difference between two correlated variables**

If we know that the variance of  $Y_1$  is  $V_1$  and the variance of  $Y_2$  is  $V_2$ , then

$$\text{var}(Y_1 - Y_2) = V_1 + V_2 - 2r\sqrt{V_1}\sqrt{V_2},$$

where  $r$  is the correlation coefficient that describes the extent to which  $Y_1$  and  $Y_2$  co-vary. If  $Y_1$  and  $Y_2$  are inextricably linked (so that a change in one determines completely the change in the other), then  $r = 1$ , and the variance of the difference is close to zero. At the other extreme, if  $Y_1$  and  $Y_2$  are unrelated, then  $r = 0$  and the variance is the sum of the individual variances. If the correlation is  $r = 0.5$  then the variance is approximately the average of the two variances.

where  $r$  is the correlation between the two outcomes. If both variances  $V_{Y1}$  and  $V_{Y2}$  are equal (say to  $V$ ), then (29.2) simplifies to

$$V_{\bar{Y}} = \frac{1}{2}V(1+r). \quad (29.3)$$

In the running example, in study 1 the effect sizes for math and reading are 0.30 and 0.10, the variance for each is 0.02. Suppose we know that the correlation between them is 0.50. The composite score for *Basic skills* ( $\bar{Y}$ ) is computed as

$$\bar{Y} = \frac{1}{2}(0.30 + 0.10) = 0.2000,$$

with variance (based on (29.2))

$$V_{\bar{Y}} = \frac{1}{4}(0.05 + 0.05 + 2 \times 0.50 \times \sqrt{0.05} \times \sqrt{0.05}) = 0.0375,$$

or, equivalently (using (29.3)),

$$V_{\bar{Y}} = \frac{1}{2} \times 0.05 \times (1 + 0.50) = 0.0375.$$

Using this formula we can see that if the correlation between outcomes was zero, the variance of the composite would be 0.025 (which is half as large as either outcome alone) because the second outcome provides entirely independent information. If the correlation was 1.0 the variance of the composite would be 0.050 (the same as either outcome alone) because all information provided by the second outcome is redundant. In our example, where the correlation is 0.50 (*some* of the information is redundant) the variance of the composite falls between these extremes. When we were working with independent subgroups (earlier in this chapter) the correlation was zero, and therefore the variance of the composite was 0.025.

These formulas are used to create Table 29.2, where the variance for each composite is based on formula (29.2) and the weight is simply the reciprocal of the variance.

At this point we can proceed to the meta-analysis using these five (synthetic) scores. To compute a summary effect and other statistics using the fixed-effect model, we apply the formulas starting with (11.3). Using values from the line labeled *Sum* in Table 29.3,

$$M = \frac{50.118}{275.542} = 0.1819,$$

**Table 29.3** Multiple outcomes – summary effect.

Study	Outcome	ES	Variance	ES	Correlation	Variance	Weight	ES*WT
Study 1	Reading	0.300	0.050	0.200	0.500	0.038	26.667	5.333
	Math	0.100	0.050					
Study 2	Reading	0.200	0.020	0.150	0.600	0.016	62.500	9.375
	Math	0.100	0.020					
Study 3	Reading	0.400	0.050	0.300	0.600	0.040	25.000	7.500
	Math	0.200	0.050					
Study 4	Reading	0.200	0.010	0.150	0.400	0.007	142.857	21.429
	Math	0.100	0.010					
Study 5	Reading	0.400	0.060	0.350	0.800	0.054	18.519	6.481
	Math	0.300	0.060					
<b>Sum</b>							<b>275.542</b>	<b>50.118</b>

with variance

$$V_M = \frac{1}{275.542} = 0.0036.$$

The average difference between the tutored and control groups on *Basic skills* is 0.1819 with variance 0.0036 and standard error 0.060. The 95% confidence interval for the average effect is 0.064 to 0.300. The Z-value for a test of the null hypothesis is 3.019 with a two-sided *p*-value of 0.003.

### **Working with more than two outcomes per study**

These formulas can be extended to accommodate any number of outcomes. If *m* represents the number of outcomes within a study, then the composite effect size for that study would be computed as

$$\bar{Y} = \frac{1}{m} \left( \sum_j^m Y_j \right), \quad (29.4)$$

and the variance of the composite is given by

$$V_{\bar{Y}} = \left( \frac{1}{m} \right)^2 \text{var} \left( \sum_{j=1}^m Y_i \right) = \left( \frac{1}{m} \right)^2 \left( \sum_{j=1}^m V_i + \sum_{j \neq k} (r_{jk} \sqrt{V_j} \sqrt{V_k}) \right) \quad (29.5)$$

as derived in Box 29.1. If the variances are all equal to *V* and the correlations are all equal to *r*, then (29.5) simplifies to

$$V_{\bar{Y}} = \frac{1}{m} V (1 + (m - 1)r). \quad (29.6)$$

### **Impact of the correlations on the combined effect**

One issue to consider is what happens as the correlation moves toward 1.0. Continuing with the simplified situation where all observations within a study have the same variance (*V*) and all pairs of observations within the study have the same correlation (*r*), if the *m* observations are independent of each other (*r* = 0), the variance of the composite is *V/m*. If the *m* observations are not independent of each other, then the variance of the composite is *V/m* times a correction factor. We will refer to this correction factor as the *variance inflation factor* (*VIF*), which is

$$VIF = 1 + (m - 1)r, \quad (29.7)$$

where *m* is the number of observations and *r* is the correlation between each pair. An increase in either *m* or *r* (or both) will result in a higher inflation of the variance compared with treating the different outcomes as independent of each other.

In Table 29.4 we explore how the variance inflation factor depends on the value of *r*, the correlation coefficient. For the purposes of this illustration we assume the simplistic situation of a study having just two outcomes (*m* = 2) with the same variance (*V* = 0.2) for each outcome. Each column in the table (A–E) corresponds to a different correlation coefficient between these outcomes.

The variance of the composite for the study is

$$V_{\bar{Y}} = \frac{1}{m} V \times VIF. \quad (29.8)$$

**Table 29.4** Multiple outcomes – impact of correlation on variance of summary effect.

	A	B	C	D	E
Effect size (here assumed identical for all outcomes)	0.4	0.4	0.4	0.4	0.4
Number of outcomes ( $m$ )	2	2	2	2	2
Variance of each outcome ( $V$ )	0.2	0.2	0.2	0.2	0.2
Correlation among outcomes ( $r$ )	0.000	0.250	0.500	0.750	1.000
Variance inflation factor ( $VIF$ )	1.000	1.250	1.500	1.750	2.000
Variance of composite	0.100	0.125	0.150	0.175	0.200
Standard error of composite	0.316	0.354	0.387	0.418	0.447
Standard error inflation factor	1.000	1.118	1.225	1.323	1.414
Lower limit	-0.220	-0.293	-0.359	-0.420	-0.477
Upper limit	1.020	1.093	1.159	1.220	1.277
p-value (2-tailed)	0.206	0.258	0.302	0.339	0.371

Taking as an example column C, the correlation is  $r = 0.50$ , and the variance inflation factor is

$$VIF = 1 + (2 - 1) \times 0.50 = 1.5000.$$

Thus, the variance of the composite score for the study is

$$V_{\bar{Y}} = \frac{1}{2} \times 0.2 \times 1.5000 = 0.150.$$

As we move from left to right in the table (from a correlation of 0.00 to 1.00) the variance inflation factor ( $VIF$ ) and (by definition) the variance double. If the inflation factor for the variance moves from 1.00 to 2.00, it follows that the inflation factor for the standard error (which is the square root of the variance) will move from 1.00 to 1.44. Therefore, the width of the confidence interval will increase by a factor of 1.44 (and, correspondingly, the Z-value for the test of the null hypothesis for this study would decrease by a factor of 1.44).

### **When the correlation is unknown**

This table also provides a mechanism for working with synthetic variables when we don't know the correlation among outcomes. Earlier, we assumed that the correlation between math and reading was known to be 0.50, and used that value to compute the standard error of the combined effect and the related statistics. In those cases where we don't know the correlation for the study in question, we should still be able to use other studies in the same field to identify a plausible range for the correlation. We could then perform a sensitivity analysis and might assume, for example, that if the correlation falls in the range of 0.50 to 0.75 then the standard error probably falls in the range of 0.39 to 0.42 (columns C to D in the table).

Researchers who do not know the correlation between outcomes sometimes fall back on either of two 'default' positions. Some will include both math and verbal

scores in the analysis and treat them as independent. Others would use the average of the reading variance and the math variance. It is instructive, therefore, to consider the practical impact of these choices.

Treating the two outcomes as independent of each other yields the same precision as setting the correlation at 0.00 (column A). By contrast, using the average of the two variances yields the same precision as setting the correlation at 1.00 (column E). In effect, then, researchers who adopt either of these positions as a way of bypassing the need to specify a correlation, are actually adopting a correlation, albeit implicitly. And, the correlation that they adopt falls at either extreme of the possible range (either zero or 1.0). The first approach is almost certain to underestimate the variance and overestimate the precision. The second approach is almost certain to overestimate the variance and underestimate the precision. In this context, the idea of working with a *plausible* range of correlations rather than the *possible* range offers some clear advantages.

As we noted at the outset, exactly the same approach applies to studies with multiple outcomes and to studies with multiple time-points. However, there could be a distinction between the two when it comes to deciding what is a plausible range of correlations. When we are working with different outcomes at a single point in time, the plausible range of correlations will depend on the similarity of the outcomes. When we are working with the same outcome at multiple time-points, the plausible range of correlations will depend on such factors as the time elapsed between assessments and the stability of the relative scores over this time period.

One issue to consider is what happens if the correlations between multiple outcomes are higher in some studies than in others. This variation will affect the *relative weights* assigned to different studies, with *more weight* going to the study with a *lower correlation*. In the running example the variances for reading and math were the same in studies 1 and 3, but the correlation between reading and math was higher in study 3. Therefore, study 3 had a higher variance and was assigned less weight in the meta-analysis.

## COMPARING OUTCOMES OR TIME-POINTS WITHIN A STUDY

We now turn to the problem of investigating *differences* between outcomes or between time-points. To extend the current example, suppose that each study reports the impact of the intervention for math and for reading and we want to know if the impact is stronger for one of these outcomes than for the other. Or, each study reports the effect at 6 months and 12 months, and we want to know if the effect changes over time.

When our goal was to compute a combined effect based on both outcomes our approach was to create a synthetic variable for each study (defined as the mean of the effect sizes) and to use this as the effect size in the analysis. We will follow the same approach here, except that the synthetic variable will be defined as the difference in effect sizes rather than their mean.

The approach is shown in Table 29.5. As before, we start with summary data for two outcomes (math and reading), and compute an effect size and variance for each.

**Table 29.5** Creating a synthetic variable as the difference between two outcomes.

Study	Outcome	Effect size	Variance	Mean	Variance
Study 1	Math	0.30	0.05	0.20	?
	Reading	0.10	0.05		

Then, we compute a synthetic effect size, which is now the *difference* between the two effects and its variance, as explained below.

This approach allows us to address the problem of correlated error, since the formula for the variance of the synthetic variable will take into account the correlation between the outcomes.

### Computing a variance for correlated outcomes

Whenever we use sample data to estimate a difference, the variance reflects the error of our estimate. If we compute the difference of two *unrelated* outcomes, each with variance  $V$ , then the variance of the difference is  $2V$ , which incorporates the two sources of error. By contrast, if we compute the difference of two (positively) related outcomes, then some of the error is redundant, and so the total error is less than  $2V$ . If the correlation between outcomes is 0.50, the variance of the difference would be equal to  $V$ , and as the correlation approaches 1.00, the variance of the difference would approach zero. The operating principle is that the higher the correlation between the outcomes, the lower the variance (the higher the precision) of the difference. The formula for the variance of a difference from correlated outcomes (see Box 29.1) is

$$V_{Y_{\text{diff}}} = V_{Y_1} + V_{Y_2} - 2r\sqrt{V_{Y_1}\sqrt{V_{Y_2}}}. \quad (29.9)$$

In words, we sum the two variances and then *subtract* a value that reflects the correlated error.

Note the difference from the formula for variance of a mean, where we *added* the correlated error and included a scaling factor. When we *combine* positively correlated outcomes, a higher correlation between outcomes results in a *higher variance*. By contrast, when we compute the *difference* between positively correlated outcomes, a higher correlation between outcomes results in a *lower variance*.

To understand why, suppose that we assess patients using two different measures of depression. If a patient is having a particularly good day when the measures are taken, both scores will tend to be higher than the patient's average. If the patient is having a bad day, both will tend to be lower than the patient's average. If we compute a *combined* effect, the error will *build* as we increase the number of outcomes. If this is a *good* day, both measures will over-estimate the patient's level of functioning. By contrast, if we compute a *difference*, we *subtract* one effect from the other, and the day-to-day variation is removed.

***Computing a difference between outcomes***

With this as background, we can return to the running example, and discuss the computation of the synthetic effect size and its variance.

The difference between reading and math is computed as

$$Y_{\text{diff}} = Y_1 - Y_2, \quad (29.10)$$

with variance

$$V_{Y_{\text{diff}}} = V_{Y_1} + V_{Y_2} - 2r\sqrt{V_{Y_1}}\sqrt{V_{Y_2}}. \quad (29.11)$$

If both variances are equal to the same variance,  $V$ , then (29.11) simplifies to

$$V = 2V(1 - r). \quad (29.12)$$

In the running example, the difference between reading and math in study 1 is computed as

$$V_{\text{diff}} = 0:30 - 0:10 = 0:2000;$$

with variance

$$V_{Y_{\text{diff}}} = 0.05 + 0.05 - 2 \times 0.50 \times \sqrt{0.05}\sqrt{0.05} = 0.05$$

or equivalently,

$$V_{Y_{\text{diff}}} = 2(0.05)(1 - 0.50) = 0.05$$

These formulas are used to create Table 29.6, where the variance for each composite is based on formula (29.11) and the weight is simply the reciprocal of the variance. At this point we can proceed to the meta-analysis using these five (synthetic) scores. The scores happen to represent difference scores, but the same formulas apply. Under the fixed-effect model the formulas starting with (11.3) yield a summary effect

$$M = \frac{27.750}{232.500} = 0.1194,$$

**Table 29.6** Multiple outcomes – difference between outcomes.

Study	Outcome	ES	Variance	ES	Correlation	Variance	Weight	ES*WT
Study 1	Reading	0.300	0.050	0.200	0.500	0.050	20.000	4.000
	Math	0.100	0.050					
Study 2	Reading	0.200	0.020	0.100	0.600	0.016	62.500	6.250
	Math	0.100	0.020					
Study 3	Reading	0.400	0.050	0.200	0.600	0.040	25.000	5.000
	Math	0.200	0.050					
Study 4	Reading	0.200	0.010	0.100	0.400	0.012	83.333	8.333
	Math	0.100	0.010					
Study 5	Reading	0.400	0.060	0.100	0.800	0.024	41.667	4.167
	Math	0.300	0.060					
Sum							232.500	27.750

with variance

$$V_M = \frac{1}{232.500} = 0.0043.$$

The average difference between the effect size for reading and the effect size for math is 0.1194 with variance 0.0043 and standard error 0.066. The 95% confidence interval for the average difference is -0.009 to 0.248. The Z-value for a test of the null hypothesis is 1.820 with a two-sided *p*-value of 0.069.

### **Working with more than two outcomes per study**

The formulas presented for a difference based on two outcomes can be extended to accommodate any number of outcomes using contrasts. For example, we could look at the difference between (a) math scores and (b) the mean of reading and verbal scores. However, this is beyond the scope of this volume.

### **Impact of the correlations on the combined effect**

One issue to consider is what happens if the correlation between two outcomes is higher in some studies than in others. This variation will affect the *relative weights* assigned to different studies, with more weight going to the study with a higher correlation. In the running example the variances for reading and math were the same in studies 1 and 3, but the correlation between reading and math was higher in study 3. Therefore, study 3 had a lower variance and was assigned *more* weight in the meta-analysis. This is the opposite of what happens for a composite.

A second issue to consider is what happens as the set of correlations as a whole moves toward 1.0. Continuing with the simplified situation where both observations within a study have the same variance ( $V$ ), if the two observations are independent of each other, the variance of the composite is  $2V$ . If the observations are not independent of each other, then the variance of the composite is  $2V$  times a correction factor. We will refer to this correction factor as the *variance inflation factor (VIF)*,

$$VIF = 1 - r, \quad (29.13)$$

where  $r$  is the correlation between the two components. An increase in  $r$  will result in a deflation of the variance compared with treating the different outcomes as independent of each other.

In Table 29.7 we explore how the variance inflation factor depends on the value of  $r$ , the correlation coefficient. For the purposes of this illustration we assume the simplistic situation of a study having the same variance ( $V_Y = 0.2$ ) for each outcome. Each column in the table (A–E) corresponds to a different correlation coefficient between these outcomes.

The variance of the difference between the outcomes is

$$V_{Y_{diff}} = 2 \times V_Y \times VIF. \quad (29.14)$$

Taking as an example column C, the correlation is  $r = 0.50$  and the variance inflation factor is

$$VIF = 1 - 0.50 = 0.5000.$$

**Table 29.7** Multiple outcomes – Impact of correlation on the variance of difference.

	A	B	C	D	E
Difference ( $\bar{Y}$ )	0.4	0.4	0.4	0.4	0.4
Variance of each outcome ( $V$ )	0.2	0.2	0.2	0.2	0.2
Correlation between outcomes	0.000	0.250	0.500	0.750	1.000
Variance inflation factor	1.000	0.750	0.500	0.250	0.000
Variance of difference ( $V_{\bar{Y}_{diff}}$ )	0.400	0.300	0.200	0.100	0.000
Standard error of difference	0.283	0.245	0.200	0.141	0.003
Standard error inflation factor	1.000	0.866	0.707	0.500	0.010
Lower limit	0.155	-0.0802	0.008	0.124	0.394
Upper limit	0.955	0.8802	0.792	0.676	0.406
p-value (2-tailed)	0.158	0.103	0.046	0.005	0.000

Thus the variance of the composite score for one study is

$$V_{\bar{Y}_{diff}} = 2 \times 0.20 \times 0.50 = 0.20.$$

As we move from left to right in the table (from a correlation of 0.00 to 0.75) the variance inflation factor (*VIF*) moves from 1.00 to 0.25. If the inflation factor for the variance moves from 1.00 to 0.25, it follows that the inflation factor for the standard error (which is the square root of the variance) will move from 1.00 to 0.50.

Therefore, the confidence interval will narrow by 50 %, and the Z-value for the test of the null hypothesis will double.

Note. In this example we focused on correlations in the range of 0.0 to 0.75, columns A–D in the table. As the correlation approaches 1.0 (column E) the variance will approach zero. This means that the width of the confidence interval will approach zero, the Z-value will approach infinity, and the p-value will approach zero. These apparent anomalies reflect what would happen if all error were removed from the equation. In Table 29.7, the correlation displayed as 1.000 is actually entered as 0.9999.

### When the correlation is unknown

Table 29.7 also provides a mechanism for working with synthetic variables when we don't know the correlation among outcomes. Earlier, we assumed that the correlation between math and reading was known to be 0.50, and used that value to compute the standard error of the difference and related statistics. In those cases where we don't know the correlation for the study in question, we should still be able to use other studies in the same field to identify a plausible range for the correlation. We could then perform a sensitivity analysis and say, for example, that if the correlation falls in the range of 0.50 to 0.75 then the two-tailed p-value probably falls in the range of 0.046 to 0.005 (columns C to D in the table).

Researchers who do not know the correlation between outcomes sometimes treat the outcomes as coming from independent subgroups. It is instructive, therefore, to consider the practical impact of this choice. Treating the two outcomes as independent

of each other yields the same precision as setting the correlation at 0.00 (column A). In effect, then, researchers who take this approach as a way of bypassing the need to specify a correlation, are actually adopting a correlation, albeit implicitly. And, the correlation that they adopt is zero. As such, it is almost certain to overestimate the variance and underestimate the precision of the difference. In this context, (as in the case of a combined effect) the idea of working with a *plausible range of correlations* offers some clear advantages.

### SUMMARY POINTS

- When we have effect sizes for more than one outcome (or time-point) within a study, based on the same participants, the information for the different effects is not independent and we need to take account of this in the analysis.
- To compute a summary effect using multiple outcomes we create a synthetic effect size for each study, defined as the mean effect size in that study, with a variance that takes account of the correlation among the different outcomes. We then use this effect size and variance to compute a summary effect across studies. Higher correlations yield *less* precise estimates of the summary effect.
- To compute the difference in effects we create a synthetic effect size for each study, defined as the *difference* between effect sizes in that study, with a variance that takes account of the correlation among the different outcomes. We then use this effect size and variance to assess the difference in effect sizes. Higher correlations yield *more* precise estimates of the difference in effects.

### Further Reading

- Cooper, H. (1982). Scientific Guidelines for conducting integrative research reviews. *Review of Educational Research* 52(2): 291–302.
- Hedges, L.V. (2019). Stochastically dependent effect sizes. In Cooper, H., Hedges, L.V. & Valentine, J.C. (Eds.) *The Handbook of Research Synthesis* (Third Edition). New York: Russell Sage Foundation.



# Multiple Comparisons within a Study

---

### Introduction

Combining across multiple comparisons within a study

Differences between treatments

---

## INTRODUCTION

The final case of a complex data structure is the case where studies use a single control group and several treatment groups. For example, suppose we are working with five studies that assessed the impact of tutoring on student performance. Each study included three groups – a control group (a free study period), intervention A (tutoring focused on that day's school lesson) and intervention B (tutoring based on a separate agenda).

If our goal was to compute a summary effect for *A* versus control and separately for *B* versus control, we would simply perform two separate meta-analyses, one using the *A* versus control comparison from each study, and one using the *B* versus control comparison from each study.

The issues we address in this chapter are how to proceed when we want to incorporate both treatment groups in the same analysis. Specifically,

- we want to compute a summary effect for the active intervention (combining *A* and *B*) versus control;
- or, we want to investigate the difference in effect size for intervention *A* versus intervention *B*.

## COMBINING ACROSS MULTIPLE COMPARISONS WITHIN A STUDY

The issue we need to address is that the effect for *A* versus control and the effect for *B* versus control are not independent of each other. If each group (*A*, *B* and control) has 200 participants and we treated the two effects as independent, our effective sample

size would appear to be 800 (since we count the control group twice) when in fact the true sample size is 600.

The problem, and the solution, are very similar to the ones we discussed for multiple outcomes (or time-points) within a study. If our goal is to compare *any treatment* versus control, we can create a composite variable which is simply the mean of *A* versus control and *B* versus control. The variance of this composite would be computed based on the variance of each effect size as well as the correlation between the two effects. At that point, all the formulas for combining data from multiple outcomes would apply here as well.

The difference between multiple outcomes and multiple comparison groups is the following. In the case of multiple outcomes, the correlation between outcomes could fall anywhere in the range of zero (or even a negative correlation) to 1.0. We suggested that the researcher work with a range of plausible correlations, but even this approach would typically yield a nontrivial range of possible correlations (say, 0.25 to 0.75) and variances. By contrast, in the case of multiple comparison groups, the correlation can be estimated accurately based on the number of cases in each group. For example, if group *A*, group *B* and the control group each have 200 participants, the correlation between *A* versus control and *B* versus control is 0.50. (This follows from the fact that the correlation between group *A* and group *B* is 0, while the correlation between control and control is 1, yielding a combined correlation midway between the two, or 0.50.) Therefore, we can work with a correlation of 0.50 without the need to conduct a sensitivity analysis based on a range of possible correlations.

This approach can be extended for the case where the sample size differs from one group to the next. For example, we can use a weighted mean of the effects rather than a simple mean, to give more weight to the treatments with more subjects. In this case, we would also need to adjust the variance to take account of the weighting. If the sample size differs from group to group the correlation will no longer be 0.50, but can be estimated precisely based on the data, without the need to resort to external correlations.

An alternate approach for working with multiple comparison groups is to collapse data from the treatment groups and use this data to compute an effect size and variance. If we are working with binary data from one control group and two treatment groups in a  $2 \times 3$  table we would collapse the two treatment groups to create a  $2 \times 2$  table, and then compute the effect size from that. Or, if we are working with means and standard deviations for three groups (*A*, *B* and control) we would collapse the data from *A* and *B* to yield a combined mean and standard deviation, and then compute the effect size for the control group versus this merged group (see option 2 for independent subgroups, (23.1), (23.2), and (23.3)). This approach will yield essentially the same results as the method proposed above.

## DIFFERENCES BETWEEN TREATMENTS

We now turn to the problem of investigating *differences* between treatments when the different treatments use the same comparison group. To extend the current example,

suppose that we have two treatment groups (*A* and *B*) and a control group. We want to know if one of the treatments is superior to the other.

While the approach used for computing a *combined* effect with multiple comparisons was similar to that for multiple outcomes, this is not the case when we turn to *differences* among the treatments. In the case of multiple outcomes, we had an effect size for reading (defined as the difference between treated and control) and an effect size for math (the difference between treated and control). Our approach was to work with the difference between the two effect sizes. In the case of multiple comparisons the analogous approach would be to compute the effect size for *A* versus control, and the effect size for *B* versus control, and then work with the difference between the two effect sizes.

While this approach would work, a better (potentially more powerful) approach is to ignore the control group entirely, and simply define the effect size as the difference between *A* and *B*. If we are working with binary data we would create a  $2 \times 2$  table for *A* versus *B* and use this to compute an effect size. If we are working with means and standard deviations we would compute an effect size from the summary data in these two groups, ignoring the control group entirely.

This approach will only work if all studies have the same groups (here, *A* and *B*), which allows us to create the same effect size (*A* versus *B*) for each study. In practice, we are likely to encounter problems since some studies might compare *A* versus *B*, while others compare *A* versus control and still others compare *B* versus control. Or, some studies might include more than two comparison groups. Methods developed to address these kinds of issues are beyond the scope of this book, but are covered in the further readings.

### SUMMARY POINTS

- When a study uses one control group and more than one treatment group, the data from the control group is used to compute more than one effect size. Therefore, the information for these effect sizes is not independent and we need to take this into account when computing the variance.
- To compute a combined effect, ignoring differences among the treatments, we can create a synthetic effect size for each study, defined as the mean effect size in that study (say, the mean of treatment *A* versus control and of treatment *B* versus control), with a variance that takes account of the correlation among the different treatments. We can then use this synthetic effect size and variance to compute a summary effect across studies. This is the same approach used with multiple outcomes.
- To look at differences among treatments the preferred option is to perform a direct comparison of treatment *A* versus treatment *B*, removing the control group from the analysis entirely. In some cases this will not be possible for practical reasons. In this case we can revert to the synthetic effect size, or can apply advanced methods.

## Further Reading

- Caldwell, D.M., Ades A.E. & Higgins, J.P.T. (2005). Simultaneous comparison of multiple treatments: combining direct and indirect evidence. *BMJ* 331: 897–900.
- Glass, G.V. McGaw, B., & Smith, M.L., (1981). *Meta-analysis in Social Research*. Beverly Hills: Sage Publications.
- Higgins J.P.T., Thomas, J., Chandler, J., Cumpston, M., Li, T., Page, M.J., Welch, V.A. (editors) (2019). *Cochrane Handbook for Systematic Reviews of Interventions*. 2nd Edition. Chichester, UK: John Wiley & Sons, Ltd.
- Salanti G, Higgins J.P.T., Ades A.E., Ioannidis J.P.A. (2008). Evaluation of networks of randomized trials. *Statistical Methods in Medical Research* 17: 279–301.

# Notes on Complex Data Structures

---

Introduction  
Summary effect  
Differences in effect

---

## INTRODUCTION

In this Part we discussed three cases where studies provide more than one unit of data for the analysis. These are the case of multiple independent subgroups within a study, multiple outcomes or time-points based on the same subjects, and two or more treatment groups that use the same comparison group.

## SUMMARY EFFECT

One issue we addressed was how to compute a summary effect using all of the data. For independent subgroups this meant looking at the impact of treatment versus no treatment, and ignoring any differences between stage-1 patients and stage-2 patients. For multiple outcomes this meant looking at the impact of the intervention on basic skills, and ignoring any differences between the impact on math versus reading. For multiple treatment groups it meant looking at the impact of treatment and control, and ignoring any differences among variants of the treatment.

In all cases, the key issue was the need to address any possible redundancy of information, since the precision of the combined effect is strongly dependent on the amount of information. To highlight the difference among the three cases, we used the same numbers in the worked examples as we moved from one chapter to the next. For example, for study 1, we assumed a variance of 0.05 between the two units of information, whether for two independent groups or for two outcomes.

The formula for the variance of a composite is the same in all cases, namely (in its simple form)

$$V_{\bar{Y}} = \frac{1}{2} V_Y (1 + r). \quad (31.1)$$

Consider, then, how this formula plays out for independent subgroups, for multiple outcomes, and for multiple comparisons.

- For multiple independent subgroups  $r$  is always zero, and so the variance of the composite is 0.025.
- For multiple outcomes  $r$  can take on any value. We assume that it falls somewhere in the range of 0.00 to 1.00, and assume further that the researcher can provide a range of plausible values for  $r$ . For example, we might have evidence that  $r$  is probably close to 0.50, with a plausible range of 0.25 to 0.75. In that case the variance would probably be close to 0.038, with a plausible range of 0.031 to 0.044.
- For multiple comparisons  $r$  can be determined based on the number of treatment groups and the number of subjects in each group. If there is one control group and two treatment groups, and the subjects are divided evenly across groups, then  $r$  would be 0.50 and the variance would be 0.038. In other cases  $r$  could move either up or down, but can always be computed, and therefore we don't need to work with a range of values.

## DIFFERENCES IN EFFECT

The second issue we addressed was how to look at the difference between effects.

For independent subgroups this meant looking at whether the treatment was more effective for one of the subgroups (stage-1 or stage-2 patients) than the other. For multiple outcomes this meant looking at whether the intervention had more of an impact on one outcome (reading or math) than the other. For multiple treatment groups it meant looking at whether one of the treatments was more effective than the other.

Again, the key issue was the need to address any possible redundancy of information, since the precision of the difference is strongly dependent on the amount of information. To highlight the difference among the three cases we used the same numbers in the worked examples as we moved from one chapter to the next. For example, we assumed a variance of 0.05 for each subgroup, or for each outcome, or for each comparison.

The formula for the variance of a difference is the same in all cases, namely (in its simple form)

$$V_{Y_{\text{diff}}} = 2V_Y(1 - r). \quad (31.2)$$

This formula incorporates the term  $(1 - r)$ , which means that a higher correlation will yield a more precise estimate of the difference. This is the reverse of the case for a composite effect, where the formula incorporates the term  $(1 + r)$ , and a higher correlation will yield a less precise estimate.

- For multiple independent subgroups  $r$  is always zero, and so the variance of the difference is 0.100.
- For multiple outcomes  $r$  can take on any value. We assume that it falls somewhere in the range of 0.00 to 1.00, and assume further that the researcher can provide a range of plausible values for  $r$ . For example, we might have evidence that  $r$  is probably

close to 0.50, with a plausible range of 0.25 to 0.75. In that case the variance would probably be close to 0.050, with a plausible range of 0.075 to 0.025.

- For multiple comparisons we have essentially the same situation as for multiple outcomes, except that we can actually compute the correlation needed for the formula. However, there are other approaches available that allow for head-to-head comparisons of the treatment groups.



## **Other Issues**



# Overview

The first chapter in this Part addresses the issue of vote counting. Vote counting is the name used to describe the idea of seeing how many studies yielded a significant result, and how many did not. We explain why this approach is always a bad idea. In fact, if the techniques used in meta-analysis are extensions of procedures in primary studies, then vote counting for multiple studies is an extension of a mistake that is ubiquitous in primary studies.

In the next chapter we address the issue of statistical power in meta-analysis. A meta-analysis often yields a more powerful test of the null hypothesis than any of the separate studies. Here, we explain why this is true, and offer some examples to show how important this can be. At the same time, we caution that this is not always the case. While there is a general perception that all meta-analyses have high power to yield a statistically significant effect, this is not always the case for the main effect, and is rarely the case for other tests (such as tests for heterogeneity). We discuss the factors that drive power in a meta-analysis, and compare these with the factors that drive power in a primary study.

Another chapter outlines the issue of publication bias. Several lines of evidence demonstrate that studies that yield larger effect sizes are more likely to be published, and incorporated in a meta-analysis, than similar studies that yield smaller effect sizes. We discuss the evidence for this phenomenon, and methods to assess its likely impact on any given meta-analysis.



# Vote Counting – A New Name for an Old Problem

---

### Introduction

Why vote counting is wrong

Vote counting is a pervasive problem

---

## INTRODUCTION

One question we often ask of the data is whether or not it allows us to reject the null hypothesis of no effect. Researchers who address this question using a narrative review need to synthesize the  $p$ -values reported by the separate studies. Since these are discrete pieces of information and the narrative review provides no statistical mechanism for synthesizing these values, narrative reviewers often resort to a process called vote counting. Under this process the reviewer counts the number of statistically significant studies and compares this with the number of statistically nonsignificant studies.

In some cases this process has been formalized, such that one actually counts the number of significant and nonsignificant  $p$ -values and picks the winner. In some variants, the reviewer would look for a clear majority rather than a simple majority. Or, the reviewer might not work directly with the  $p$ -values, but with the discussion section of the papers which are based on the  $p$ -values.

One might think that summarizing  $p$ -values through a vote-counting procedure would yield more accurate decision than any one of the single significance tests being summarized. This is not generally the case, however. In fact, Hedges and Olkin (1980) showed that the power of vote-counting considered as a statistical decision procedure can not only be lower than that of the studies on which it is based, the power of vote counting can tend toward zero as the number of studies increases. In other words, vote counting is not only misleading, it tends to be *more* misleading as the amount of evidence (the number of studies) increases!

In any event, the idea of vote counting is fundamentally flawed and the variants on this process are equally flawed (and perhaps even more dangerous, since the basic

flaw is less obvious when hidden behind a more complicated algorithm or is one step removed from the  $p$ -value). Our goal in this chapter is to explain why this is so, and to provide a few examples.

## WHY VOTE COUNTING IS WRONG

The logic of vote counting says that a significant finding is evidence that an effect exists, while a nonsignificant finding is evidence that an effect is absent. While the first statement is true, the second is not. While a nonsignificant finding *could* be due to the fact that the true effect is nil, it can also be due simply to low statistical power.

Put simply, the  $p$ -value reported for any study is a function of the observed effect size and the sample size. Even if the observed effect is substantial, the  $p$ -value will not be significant unless the sample size is adequate. In other words, as most of us learned in our first statistics course, *the absence of a statistically significant effect is not evidence that an effect is absent*.

For example, suppose five randomized controlled trials (RCTs) had been performed to test the impact of an intervention, and that none were statistically significant (the  $p$ -value in each case is 0.265) as illustrated in Figure 33.1. The vote count is 5 to 0 against an effect, and one might assume that the intervention has no effect.

By contrast, the meta-analysis (Figure 33.1), by combining the information into a single analysis, allows us to perform a proper test of the null hypothesis. Not only is this approach valid, but the test of the summary effect is often much more powerful than tests performed on any of the separate studies. When we merge the data, the effect size stays the same, but the confidence interval narrows and no longer includes the null value. The  $p$ -value for each study alone is 0.265, but the  $p$ -value for the summary effect is 0.013. Clearly, the absence of significance in each study is due to a lack of precision rather than a small effect.

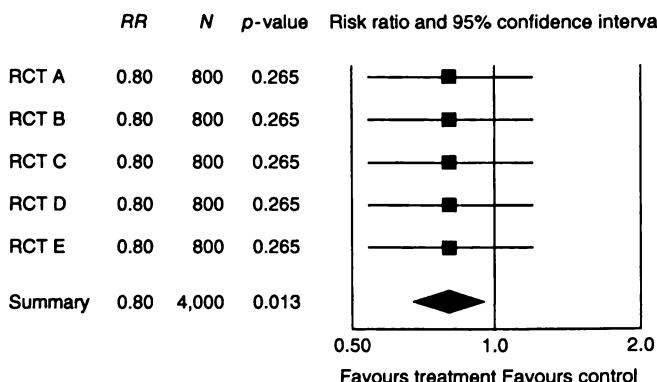


Figure 33.1 The  $p$ -value for each study is  $> 0.20$  but the  $p$ -value for the summary effect is  $< 0.02$ .

For purposes of explaining why vote counting is a bad idea, we could end the chapter here. However, because vote counting in its various forms is so pervasive, we will expand on this idea to show how the basic mistake that underlies vote counting affects much of the literature, and how meta-analysis can help address this problem.

## VOTE COUNTING IS A PERVERSIVE PROBLEM

While the term vote counting is associated with narrative reviews it can also be applied to the single study, where a significant  $p$ -value is taken as evidence that an effect exists, and a nonsignificant  $p$ -value is taken as evidence that an effect does not exist. Numerous surveys in a wide variety of substantive fields have repeatedly documented the ubiquitous nature of this mistake.

In medicine, for example, Freiman, Chalmers, Smith and Kuebler (1978) surveyed reports of controlled clinical trials that had been published in a number of medical journals (primarily *The Lancet*, the *New England Journal of Medicine*, and the *Journal of the American Medical Association* during the period 1960–1977), and selected 71 that had reported negative results. The authors found that if the true drug effect had been in the region of 50% (e.g. a mortality rate of 30% for placebo vs. 15% for drug), median power would have been 60%. In other words, even if the drug cut the mortality rate in half there was still a 40% probability that the study would have failed to obtain a statistically significant result.

The authors went on to make the following point: Despite the fact that power was terribly low, in most cases the absence of statistical significance was interpreted as meaning that the drug *was not effective*. They wrote: ‘The conclusion is inescapable that many of the therapies discarded as ineffective after inconclusive “negative” trials may still have a clinically meaningful effect’ (p. 694). In fact, it is possible (or likely) that some of the therapies discarded on this basis might well have had very substantial therapeutic effects.

In the social sciences Cohen (1962) surveyed papers published in the *Journal of Abnormal and Social Psychology* in 1960. Mean power to detect a small, medium, or large effect, respectively, was 0.18, 0.48, and 0.83. Cohen noted that despite the low power, when the studies with *negative* results are published, readers tend to interpret the absence of statistical significance as evidence that the treatment has been proven ineffective.

In the years that followed a kind of cottage industry developed of publishing papers that documented the fact of low power in any number of journals in the area of behavioral research. Many of these are cited in Sedlmeier and Gigerenzer (1989) and Rossi (1990). Similar papers were published to document the same problem in the field of medicine (Borenstein, 1994; Hartung, Cottrell & Giffen, 1983; Phillips, Scott, & Blasczynski, 1983; Reed & Slachert, 1981; Reynolds, 1980) and psychiatry (Kane & Borenstein, 1985).

Sedlmeier and Gigerenzer (1989) published a paper entitled *Do studies of statistical power have an effect on the power of statistical studies?* They found that in the

25 years since Cohen's initial survey power had not changed in any substantive way. Similarly, Rossi (1990) reviewed papers published in 1982 in the *Journals of Abnormal Psychology, Consulting and Clinical Psychology, and Personality and Social Psychology*. Mean power to detect small, medium, and large effects, respectively, was 0.17, 0.57, and 0.83.

This led one of the current authors (Borenstein, 2000) to propose four theorems, as follows.

1. Power in many fields of research is abysmally low.
2. Rule (1) appears to be impervious to change.
3. The absence of significance should be interpreted as *more information is required* but is interpreted in error as meaning *no effect exists*.
4. Rule (3) appears to be impervious to change.

In a sense, then, vote counting did not originate with the narrative review. Rather, the basic mistake has existed for decades, where it found a home in primary research. When the field moved on to narrative reviews, this basic mistake was named and codified but remained basically unchanged.

There is, however, one important difference. When we are working with a single study and we have a nonsignificant result we don't have any way of knowing whether or not the effect is real. The nonsignificant *p*-value could reflect either the fact that the true effect is nil *or* the fact that our study had low power. While we caution against accepting the former (that the true effect is nil) we cannot rule it out.

By contrast, when we use meta-analysis to synthesize the data from a series of studies we can often identify the true effect. And in many cases (for example if the true effect is substantial and is consistent across studies) we can assert that the nonsignificant *p*-value in the separate studies was due to low power rather than the absence of an effect.

In the streptokinase meta-analysis on page 10, for example, it is clear that the treatment does reduce the risk of death. It is fair to say that the reason that 27 studies had nonsignificant *p*-values was *not* because the treatment had no effect, but rather was because of low statistical power. (In the next chapter we actually compute the power for the streptokinase studies.)

### Moving beyond the null hypothesis

In this chapter we have shown that *if our goal* is to test the null hypothesis, then meta-analysis (unlike the narrative review) provides a statistically sound mechanism for this purpose. However, we want to emphasize that meta-analysis allows us *to move beyond a test of the null hypothesis*. It allows us to assess the magnitude of the effect (which is often a more relevant question) and to determine whether or not the effect size is consistent across studies.

**SUMMARY POINTS**

- Vote counting is the process of counting the number of studies that are statistically significant and comparing this with the number that are not statistically significant.
- Vote counting treats a nonsignificant  $p$ -value as evidence that an effect is absent. In fact, though, small, moderate, and even large effect sizes may yield a non-significant  $p$ -value due to inadequate statistical power. Therefore, vote counting is never a valid approach.



# Power Analysis for Meta-Analysis

---

Introduction

A conceptual approach

In context

When to use power analysis

Planning for precision rather than for power

Power analysis in primary studies

Power analysis for meta-analysis

Power analysis for a test of homogeneity

---

## INTRODUCTION

A common goal in research (both in primary studies and in meta-analysis) is to test the null hypothesis of no effect. If that is our goal, then it is important to ensure that the study has good statistical power (a sufficiently high likelihood of yielding a statistically significant result). In this chapter we pursue two themes related to statistical power.

The first theme is conceptual. We discuss the factors that determine power and explore how the value of these factors may change as we move from a primary study to a meta-analysis. On this basis we can see why power for a meta-analyses is sometimes (but not always) higher than power in any of the included studies. Here, we address only power for a test of the main effect.

The second theme is practical. We briefly review the process of power analysis for primary studies, and then show how the same process can be extended for meta-analysis. Here, we focus primarily on power for a test of the main effect, but also include material on tests of heterogeneity.

## A CONCEPTUAL APPROACH

### Background

The significance test (performed *after* the study is completed) takes the form

$$Z = \frac{M}{SE_M}, \quad (34.1)$$

where  $Z$  is the test statistic,  $M$  is the effect size, and  $SE_M$  is the standard error of the effect size. The observed value of  $Z$  is then compared with the criterion alpha ( $\alpha$ ). If alpha has been set at 0.05, then a  $p$ -value lower than 0.05 (a  $Z$ -value beyond  $\pm 1.96$ ) will be statistically significant.

Therefore, the likelihood that the results will be statistically significant depends on the following.

- The effect size. As  $M$  increases,  $Z$  increases, and the likelihood of statistical significance increases.
- The precision of the estimate. As the precision increases (as  $SE_M$  decreases),  $Z$  increases, and the likelihood of statistical significance increases.
- The criterion for significance ( $\alpha$ ). As  $\alpha$  moves away from zero, the likelihood that  $p$  will be less than  $\alpha$  increases, and the likelihood of statistical significance increases.

The difference between the significance test and a power analysis is that we perform the significance test *after* the data have been collected, at which point  $M$ ,  $SE_M$ , and  $\alpha$  are known. By contrast, when we compute power (usually, *before* the study has been performed) we need to make an assumption about  $M$ , an educated guess about  $SE_M$ , and select a value for  $\alpha$ . As such, we are working with projected values rather than observed values. Still, the factors that control power are the same as those that control the significance test. To wit, as the expected effect size increases, as the precision of the estimate increases, and/or as  $\alpha$  is moved away from zero, the higher the power.

### Power for meta-analyses as compared with primary studies

With this as background we can anticipate how power considerations differ between primary studies and meta-analysis. Typically, the expected effect size and the criterion for significance ( $\alpha$ ) will be the same for the primary study and for the meta-analysis. However, the precision will differ (sometimes substantially) as we move from primary studies to a meta-analysis.

To get a sense of how the precision changes as we move from the individual studies to the summary effect we can work with the confidence intervals on the forest plot. Concretely, the width of each confidence interval (the distance from the lower limit to the upper limit) is proportional to the standard error (the denominator in (34.1)).

For example, consider the Cannon *et al.* meta-analysis that we discussed in Chapter 1. The width of the confidence interval is substantially narrower for the summary effect than for any of the included studies. Therefore, power to test the summary effect in the meta-analysis is substantially higher than power to test the effect in any of the primary studies. Assuming a risk ratio of 85% and a baseline risk of 9.4% (the mean values actually observed in these studies), power for the four studies (Prove-it, A to Z, TNT, Ideal) was 36%, 39%, 70%, and 65%, respectively. Therefore, it is not surprising that three of the studies failed to meet the criterion for significance, with only one (TNT) yielding a  $p$ -value under 0.05. By contrast, if we synthesize the effects in a meta-analysis, the combined sample size is 13,774, the statistical power is 83%, and the observed  $p$ -value is  $< 0.0001$ .

Similarly, in Chapter 2 we discussed a meta-analysis of 33 trials that tested the impact of streptokinase (versus a placebo) to prevent death following a heart attack. Again, the confidence interval is narrower for the summary effect than for any of the included studies, and (except for the largest primary studies) the difference is substantial. Assuming that streptokinase actually reduces the risk of death by some 20% (which is the summary risk ratio for this data), then based on the sample size and event rates in the studies we can determine that only three had power exceeding 80%. Of the remaining 30 studies, none had power in the range of 60% to 80%, 3 had power in the range of 40% to 60%, 6 had power in the range of 20% to 40%, and 21 had power of less than 20%. Therefore, it is not surprising that only six of the 33 studies were statistically significant. By contrast, when we combine the effects from these studies in a meta-analysis the power for the summary effect exceeds 99.9% and the *p*-value for the summary effect is 0.0000008.

The fact that a meta-analysis will often have high power is important because (as in these examples) primary studies often suffer from low power. While researchers are encouraged to design studies with power of at least 80%, this goal is often elusive. Many studies in medicine, psychology, education and an array of other fields have power substantially lower than 80% to detect large effects, and substantially lower than 50% to detect smaller effects that are still important enough to be of theoretical or practical importance. By contrast, a meta-analysis based on multiple studies will have a higher total sample size than any of the separate studies and the increase in power can be substantial.

The problem of low power in the primary studies is especially acute when looking for adverse events. The problem here is that studies to test new drugs are *powered* to find a treatment effect for the drug, and do not have adequate power to detect side effects (which have a much lower event rate, and therefore lower power). For example, a recent meta-analysis synthesized data from 44 studies and suggested that *Avandia* may increase the risk of death or myocardial infarction (MI). Because the risk of death or MI is very low, power in the individual studies was quite low. Two of the studies had power of 15% and 18%. The other 42 studies had power of less than 8%. (These computations are based on the event rates and effect sizes combined across the 44 studies.) By contrast, the meta-analysis had power of 66% and a *p*-value that either met or approached the 0.05 criterion (this depends on the method used in the analysis).

While the examples such as the streptokinase analysis and the statins analysis, where the meta-analysis had high power, are representative of *many* meta-analyses, they are not representative of *all* meta-analyses. Assuming a nontrivial effect size, power is primarily a function of the precision, and so to understand why some meta-analyses will have high power while others (with a similar effect size) will not, we need to consider the factors that control precision, as follows.

### Power under the fixed-effect model

When we are working with a fixed-effect analysis, precision for the summary effect is always higher than it is for any of the included studies. Under the fixed-effect analysis precision is largely determined by the total sample size (accumulated over all studies in the analysis), and it follows the total sample size will be higher across studies than

within studies. If the analysis includes a large number of small studies, the difference in power can be substantial.

Consider a meta-analysis of  $k$  studies with the simplest design, such that each study comprises a single sample of  $n$  observations with standard deviation  $\sigma$ . When the effect size is *consistent* across studies we saw in Chapter 13 that the standard error of the mean,  $SE_M$ , is given by (Box 13.1)

$$SE_M = \sqrt{\frac{\sigma^2}{k \times n}}. \quad (34.2)$$

Because  $k$  (the number of studies) and  $n$  (the sample size in each study) appear together,  $SE_M$  (and therefore the power) of a fixed-effect meta-analysis will depend only on the total sample size  $k \times n$ . The power will be the same for a meta-analysis with 10 studies of 100 persons each, as it would be for one primary study with 1000 persons.

As the term  $k \times n$  approaches infinity the standard error will approach zero and (provided that the effect size is nonzero) power will approach 100%.

### Power under the random-effects model

By contrast, when we move to a random-effects analysis we need to deal with two sources of error. One is the error within studies, and the other is the variance between studies. (The latter is *real* variance, but we refer to it here as error in the sense that it leads to uncertainty in the value of the mean effect.) Now, we can no longer use sample size as a surrogate for precision. Rather, we need to compute the precision based on both components.

Under the random-effects model we saw in Chapter 13 that the standard error of the summary effect is given by

$$SE_{M^*} = \sqrt{\frac{\sigma^2}{k \times n} + \frac{\tau^2}{k}}. \quad (34.3)$$

The first term is identical to that for the fixed-effect model and, again, with a large enough sample size (either enough studies or a large enough sample within studies), this term will approach zero. By contrast, the second term (which reflects the between-studies variance) will only approach zero if the estimated value of  $\tau^2$  is zero, or as the number of *studies* approaches infinity.

These formulas do not apply exactly in practice, but the conceptual point does. Concretely, in a random-effects meta-analysis, power depends on within-study error and between-studies variation. If the effect sizes are reasonably consistent from study to study, and/or if the analysis includes a substantial number of studies, then the second of these will tend to be small, and power will be driven by the cumulative sample size. In this case the meta-analysis will tend to have higher power than any of the included studies. This was the case for the statins analysis and for the streptokinase analysis. However, if the effect size varies substantially from study to study, and the analysis includes only a few studies, then this second aspect will limit the potential power of the meta-analysis. In this case, power could be limited to some low value even if the analysis includes tens of thousands of persons.

Above, we suggested that one can get a sense for the power of the summary effect (as compared with power for the included studies) by comparing the confidence interval width in the two. The same logic extends to the difference in power between fixed-effect and random-effects analyses. If the width of the confidence interval is roughly the same in the two, then power will be roughly the same. By contrast, if the random-effects interval is substantially wider then power will be lower and, depending on the effect size, may not approach acceptable levels.

## IN CONTEXT

### Power to test main effects

There is a general perception that meta-analyses have high power to detect main effects, probably stemming from some well known meta-analyses that did include large numbers of studies and therefore had high power. The Cannon *et al.* and the streptokinase studies are two examples. A number of recent reviews in the social sciences (e.g. Wilson *et al.*, 2003a, 2003b) have included more than 200 studies each.

However, most meta-analyses have far fewer studies than this. The Cochrane Database of Systematic Reviews is a database of systematic reviews, primarily of randomized trials, for medical interventions in all areas of healthcare, and currently includes over 3000 reviews. In a survey of this database (Davey *et al.*, 2011), the median number of trials included in a meta-analysis was three. When a review includes only six studies, power to detect even a moderately large effect, let alone a small one, can be well under 80%. While the median number of studies in a review differs by the field of research, in almost any field we do find some reviews based on a small number of studies, and so we cannot simply assume that power is high.

### Power for tests comparing subgroups, and for meta-regression

Even when power to test the main effect is high, many meta-analyses are not concerned with the main effect at all, but are performed solely to assess the impact of covariates (or moderator variables). For example, we might know that an intervention increases the survival time for patients with a specific form of cancer. The question to be addressed is not whether the treatment works, but whether one variant of the treatment is more effective than another variant.

The test of a moderator variable in a meta-analysis is akin to the test of an interaction in a primary study, and both suffer from the same factors that tend to decrease power. First, the *effect size* is actually the difference between the two effect sizes and so is almost invariably smaller than the main effect size. Second, the sample size within groups is (by definition) smaller than the total sample size. Therefore, power for testing the moderator will often be very low (Hedges and Pigott, 2004). The fact of low power for tests of the effects of covariates is especially important because these kinds of analyses are often carried out in order to demonstrate that the moderator variables don't have an effect – i.e. to *accept* the null hypothesis. The logic of accepting the null hypothesis is based on the assumption of high power, and that assumption is rarely tested.

### Power for tests of homogeneity, or goodness of fit

Typically, an analysis that looks at the main effect is accompanied by a test of homogeneity. Here, a nonsignificant  $p$ -value might be taken to mean that the treatment effect is consistent across studies. Similarly, an analysis that looks at covariates (by comparing subgroups or using meta-regression) is often followed by a test for goodness of fit, and a nonsignificant  $p$ -value is taken to mean that the covariates explain all the variance. In fact, though, these kinds of analyses routinely suffer from low power. Power analysis could help researchers to recognize this fact, and refrain from drawing (possibly incorrect) conclusions based on the fact that the results are not statistically significant.

### WHEN TO USE POWER ANALYSIS

In primary studies, power analysis is used primarily to determine an appropriate sample size, since this is largely under the control of the investigators. For example, researchers might modify the sample size by changing the planned enrollment period, or the number of sites, or by changing the inclusion criteria for the study. In meta-analysis the issues are rather different because the studies already exist, but there are some parallels. We have some control over the number of studies, since the inclusion criteria could be modified to include more or fewer studies. For example, if it looks as if we will not be able to obtain enough studies to yield adequate power we might consider widening the inclusion criteria, to include studies with a wider range of types of participants. Conversely, if there appears to be an abundance of studies we may elect to narrow the inclusion criteria. Altering the inclusion criteria changes the question being addressed by the review, however, and in most situations it is more appropriate to answer the original question than to adjust it because of power considerations.

While most meta-analyses are planned after the primary studies have already been performed, there are some initiatives underway to plan the meta-analyses prospectively. For example, a consortium of hospitals may plan to perform a series of studies with the goal of incorporating all of these studies into a meta-analysis. In this case, the goal might be to ensure good power for the summary effect, rather than for the individual studies. In this situation, the meta-analyst may have direct influence over the sample sizes of the individual studies and the number of studies (or hospitals). We have seen in this chapter that power considerations depend on the intended meta-analysis model. In particular, for a fixed-effect meta-analysis only the total sample size across studies is important, whereas for a random-effects meta-analysis the option of using five hospitals with 1000 patients each (for a total of 5000 patients) versus ten hospitals with 500 patients each (for the same total) can yield substantially different values for power.

A power analysis should be performed when the review is being planned, and not after the review has been completed. Researchers sometimes conduct a power analysis after the fact, and report that *Power was low, and therefore the absence of a significant effect is not informative*. While this is correct, it is preferable to address the same

question by simply reporting the observed effect size with its confidence interval. For example, *The effect size is 0.4 with a confidence interval of -0.10 to + 0.90* is much more informative than the statement that *Power was low*. The statement of effect size with confidence intervals not only makes it clear that we cannot rule out a clinically important effect, but also gives a range for what this effect might be (here, as low as -0.10 and as high as + 0.90).

## PLANNING FOR PRECISION RATHER THAN FOR POWER

A power analysis can focus directly on power (the likelihood of a test giving a statistically significant result) or may focus on precision (the likelihood that a confidence interval will be a specific width). The two approaches are closely related, although the latter is more straightforward. As we shall see below, estimating precision is an early step on the way towards estimating power. We shall define precision formally as one divided by the variance.

## POWER ANALYSIS IN PRIMARY STUDIES

The formulas for power analysis are very similar for meta-analysis and for primary studies. Before turning to meta-analysis we review some key issues in power analysis for primary studies, to provide a context.

The statistical significance computed *after* the data are collected is a function of three elements: the true effect size, the precision with which it is estimated (strongly dependent on sample size) and the criterion used for statistical significance (alpha). A study is more likely to be statistically significant if the effect size is large, the precision is high or criterion for statistical significance is liberal (i.e. alpha is large). Power, which is simply a prediction about statistical significance, is determined by the same three elements in the same way.

This parallel is obvious in the formulas for significance and for power. Consider a significance test based on a test statistic, Z. The observed value is computed as

$$Z = \frac{M}{\sqrt{V_M}}, \quad (34.4)$$

where M is an estimate of effect size from the data and V<sub>M</sub> its variance. This Z value is evaluated with reference to the standard normal distribution, with a two tailed p-value of

$$p = 2[1 - \Phi(|Z|)]. \quad (34.5)$$

In power analysis we consider the distribution of Z, not under the null hypothesis, but under specific alternatives. We will use a parameter lambda ( $\lambda$ ) to represent an alternative true value of Z, defined as

$$\lambda = \frac{\delta}{\sqrt{V_\delta}}, \quad (34.6)$$

where  $\delta$  is the true effect size, and  $V_\delta$  its variance.

Note: we are using  $V_\delta$  for both the true variance in (the Z equation) and  $\lambda$  in (the lambda equation). Strictly, these are different and might be notated differently. However, in meta-analysis we do not usually distinguish between the two (we assume the estimated variances are known) so we will not in this chapter.

The power of a study is the probability of observing values of Z that are statistically significant when the true mean of Z is  $\lambda$ . For a two-tailed test,

$$\text{Power} = (1 - \Phi(c_\alpha - \lambda)) + \Phi(-c_\alpha - \lambda), \quad (34.7)$$

where  $c_\alpha$  is the critical value of Z associated with significance level  $\alpha$  (thus, for  $\alpha = 0.05$ ,  $c_\alpha = 1.96$ ).

The actual computation of power for a given effect size ( $\delta$ ), precision ( $1/V_\delta$ ), and  $\alpha$  is usually straightforward. In power analysis the challenge is to identify a plausible range of values for each of these factors in order to determine values for the power of the study in realistic circumstances. The same formulas are used for sample size calculations, in which the power, effect size and  $\alpha$  are fixed, and the required value of  $V_\delta$  is computed. From  $V_\delta$ , a sample size to meet these criteria can be computed. In the following discussion we focus mainly on the investigation of power rather than the computation of sample size.

### Finding a range of values for the effect size

The effect size used to compute power is defined as the true (population) effect size. However, since we don't actually know the population effect size we must select a number to serve that function. In a sample size calculation, we use the smallest effect size that it would be important to detect, and this depends entirely on the context of the study. For example, if we are looking at a treatment to reduce risk of death we might decide that even a small effect is important. By contrast, if we are looking at treatment to reduce the risk of flu, we might decide that only a relatively large effect is worth detecting.

In power analysis, we must pick plausible values for the effect size, and usually a range of values is investigated. This range should be based on substantive or clinical importance, and where possible should be informed by available data. For example, if we are working with a class of drugs that typically yields a 10–15% reduction in events, then we should probably use an effect size that falls in (or near) this range. Prior studies or a pilot study can help to identify a plausible range, and researchers may also fall back on the use of conventions, such as Cohen's proposals for small, medium and large effects in social science research (Cohen, 1987).

### Finding a range of values for the precision

The major determinant of precision is sample size. The range of possible sample sizes will be determined by practical constraints, which will vary substantially from one study to the next. For a primary experimental study, a pilot study may provide information about the potential sample size by showing how many people can be enrolled over a given period of time.

Other determinants of precision depend on the type of data being collected. For example, for continuous outcomes, precision depends on the inherent variability across individuals, expressed using the standard deviation.

### Finding a range of values for significance level ( $\alpha$ )

In any significance test we need to balance the risk of a type I error (the true effect is zero but we reject the null hypothesis) and a type II error (the true effect is not zero but we fail to reject the null hypothesis). As we move the criterion for significance (alpha) toward zero we reduce the risk of a type I error but increase risk of a type II error.

There is a common convention that alpha should be set at 0.05 (or 5%) and power at 80%, which makes sense if the damage associated with a type I error is four times as severe as that associated with a type II error (with power at 80%, the risk of a type II error is 20%, which is four times the risk of a type I error). In fact, though, it would be better to adjust these risks as appropriate for a given study. For example, if a positive result will be seen as definitive and result in an immediate change in clinical practice, we would want to protect against a type I error and would therefore use a conservative value of alpha (say, 0.01). On the other hand, if a significant result will lead to other trials to replicate the effect, while a nonsignificant result will discourage further research of the treatment, we might be more concerned with a type II error. In this case we might use a more liberal value of alpha (say, 0.10) as a mechanism to increase power.

### Illustrative example

Suppose that we plan to investigate the impact of a drug to reduce pain from migraine headaches. Patients will be allocated at random to either treatment or a placebo, and assessed after two hours on a 100-point scale. The standard deviation of scores on this scale is 10. We plan to perform this study in a hospital's migraine clinic, where we can enroll about 10 patients a month, and we wish to use a significance level of  $\alpha = 0.05$ . Suppose it is proposed that the study runs for five months. During this period we can recruit 25 patients per group. Patients report that a difference of 2 to 4 points on this scale (corresponding to a standardized mean difference of 0.2 or 0.4) would represent a meaningful effect to them, and this effect is consistent with data from previous studies.

The effect size index is  $d$  with two independent groups, and so we compute the variance using (4.20),

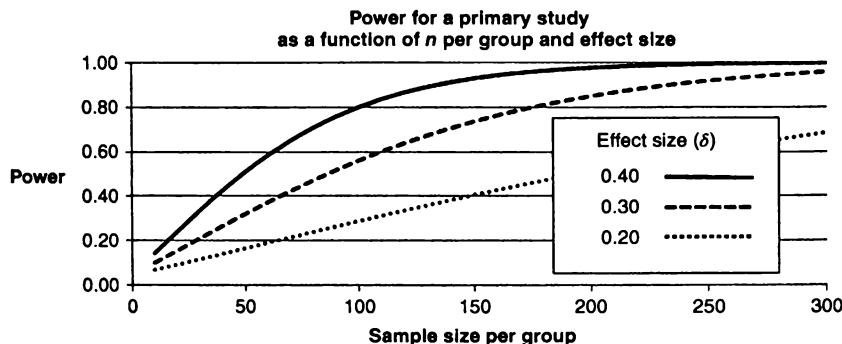
$$V_d = \frac{n_1 + n_2}{n_1 \times n_2} + \frac{d^2}{2(n_1 + n_2)}. \quad (34.8)$$

For our sample size of 25 and effect size of 0.3, we can estimate this variance as

$$V_d = \frac{25 + 25}{25 \times 25} + \frac{0.30^2}{2(25 + 25)} = 0.0809.$$

The parameter  $\lambda$  is given by

$$\lambda = \frac{0.30}{\sqrt{0.0809}} = 1.0547.$$



**Figure 34.1** Power for a primary study as a function of  $n$  and  $\delta$ .

The Z-value required for significance with alpha (2-tailed) of 0.05 is

$$c_\alpha = \Phi(1 - 0.05/2) = 1.96$$

In Excel = NORMSINV(1 - 0.05/2) = returns 1.96, and finally,

$$\text{Power} = 1 - \Phi(1.96 - 1.0547) + \Phi(-1.96 - 1.0547) = 0.1840$$

In Excel, = 1-NORMSDIST(1.96 - 1.0547) + NORMSDIST(-1.96 - 1.0547) = 0.1840.

In words, the power to detect an effect size of 0.30 is 0.184, or 18.4%. This is obviously much lower than the values of 80% or 90% that are typically desirable in a randomized trial. By applying the same formula while varying the effect size (0.20, 0.30, 0.40) and sample size (from 10 to 300 per group) we can create a graph that displays the power of the study under different scenarios.

Figure 34.1 shows power as a function of sample size for  $\delta = 0.40, 0.30$ , and  $0.20$ , assuming alpha, two-tailed, is set at 0.05. Reading left to right at 0.90 on the y-axis, to yield power of 90% for a large, medium or small effect we would need a sample size of around 140, 240, or more than 300 per group. We might then decide to enroll 240 patients per group, which will allow us to complete the study in 48 months. Reading from top to bottom at 240 on the x-axis, the study will have power of about 99% to detect the larger effect of 0.40, 90% to detect the moderate effect of 0.30, and 60% to detect the smaller effect of 0.20.

With this as prologue we can turn to meta-analysis and consider how the process is analogous to that for the primary study, and how it differs.

## POWER ANALYSIS FOR META-ANALYSIS

### Power analysis for a main effect

The logic of power analysis for meta-analysis is very similar to the logic of power analysis for primary studies. Again, we could either investigate how power is likely

to depend on plausible values of the effect size, precision and alpha, or we could compute a precision for a given effect size, alpha and power. In meta-analysis, the precision reflects both the sample sizes of the studies *and* the number of studies. When using power analysis to compute the precision, we therefore need to consider both aspects of sample size. There are differences in the implications of precision of sample sizes, depending on whether we adopt a fixed-effect or random-effects model for the meta-analysis.

For power analysis, plausible values for the effect size and precision are needed. Alpha should be chosen to reflect the potential impact of a type I error. The effect sizes should be based on substantive issues ('What effect size would it be important to detect?'). Both the effect size and the precision could be informed by a pilot study. In a primary study this might mean actually performing the study on a small scale to get a sense of the likely effect size as well as the number of persons who might be recruited. In the meta-analysis this might mean locating and coding a subset of the literature to get a sense of the effect sizes, of the sample sizes within studies, and also of the number of studies that meet the inclusion criteria.

In this chapter, we assume that every study has the same precision (in the notation of Part 3,  $V_{Y_i}$  is the same in every study). For a discussion of power analysis when the  $V_{Y_i}$  are not equal, see Hedges and Pigott (2001).

### Power for main effect using fixed-effect model

The formulas for significance and for power have the same structure as they did for primary studies. Concretely, the test of significance for the main effect in a fixed-effect meta-analysis is based on a test statistic  $Z$ , computed as

$$Z = \frac{M}{\sqrt{V_M}}, \quad (34.9)$$

but  $M$  and  $V_M$  are now the effect size and variance of the *summary effect* observed in the synthesis rather than for a single study. Recall that  $V_M$  is calculated as one divided by the sum of the weights awarded to the individual studies, where the weights are inverse-variances. If all studies had the same variance, say  $V_Y$ , then  $V_M$  is equivalent to  $V_Y/k$ , where  $k$  is the number of studies.

As before,  $Z$  is evaluated with reference to the standard normal distribution which yields the corresponding  $p$ -value. For a two-tailed test,

$$p = 2 [1 - (\Phi(|Z|))]. \quad (34.10)$$

Similarly, power analysis is again based on lambda ( $\lambda$ ), defined as

$$\lambda = \frac{\delta}{\sqrt{V_\delta}}, \quad (34.11)$$

but now  $\delta$  and  $V_\delta$  are the (hypothesized) *true* effect size and its variance for the summary effect. As before, the power formula for a two-tailed test is

$$\text{Power} = 1 - \Phi(c_\alpha - \lambda) + \Phi(-c_\alpha - \lambda). \quad (34.12)$$

### Illustrative example

Above, we presented the power analysis for a primary study to assess the impact of a treatment for migraine headaches. For a single study with  $n = 25$  per group, assuming an effect size of 0.30, we computed the variance as 0.0809, lambda as 1.0547, and power as 0.18.

Suppose that we are planning a meta-analysis to address the same question and instead of a single study with  $n = 25$  per group, we have ten studies where each has this sample size. Since the variance of a single study is

$$\frac{25 + 25}{25 \times 25} + \frac{0.30^2}{2(25 + 25)} = 0.0809,$$

the variance of the summary effect is

$$V_M = \frac{0.0809}{10} = 0.00809.$$

Then, lambda is

$$\lambda = \frac{0.30}{\sqrt{0.00809}} = 3.3354,$$

$$C_\alpha = \Phi(1 - 0.05/2) = 1.96,$$

and

$$\text{Power} = 1 - \Phi(1.96 - 3.3354) + \Phi(-1.96 - 3.3354) = 0.9155.$$

In Excel, =NORMSINV(1 - 0.05/2) returns 1.96, and=1 - NORMSDIST(1.96 - 3.3354) + NORMSDIST(-1.96 - 3.3354) returns 0.9155.

In other words, while the formula for power is identical to the formula for a primary study, the variance is reduced by a factor of  $k$  (the number of studies), which yields a  $k$ -fold increase in lambda and an increase in power.

We can apply the formula while varying the effect size and the number of studies to graph power as a function of effect size and number of studies. In Figure 34.2 we assume an  $n$  of 25 per group within each study, alpha (2-tailed) of 0.05 and a fixed-effect model. The graph shows power for an effect size ( $\delta$ ) of 0.4, 0.3, 0.2, and the number of studies varies from 1 to 25.

Reading from left to right, to ensure power of 90% for an effect size of 0.40, 0.30 or 0.20 we would need 6 studies, 10 studies, or 22 studies. Suppose that it becomes clear from a pilot study that there are likely to be at least 20 studies that meet the inclusion criteria. At that point, the researchers might decide to proceed as planned. Based on this number, and reading the graph from top to bottom, if the true effect is 0.40 or 0.30 power will exceed 99%, and if the true effect is 0.20, power will approach 90%.

By contrast, if it seemed that there would be only 5 or 10 studies that met the inclusion criteria the reviewers might wish to obtain a more accurate estimate of the number (the difference between 5 and 10 is important). Alternatively, there may be reasonable ways of modifying the research question to increase the number of relevant studies.

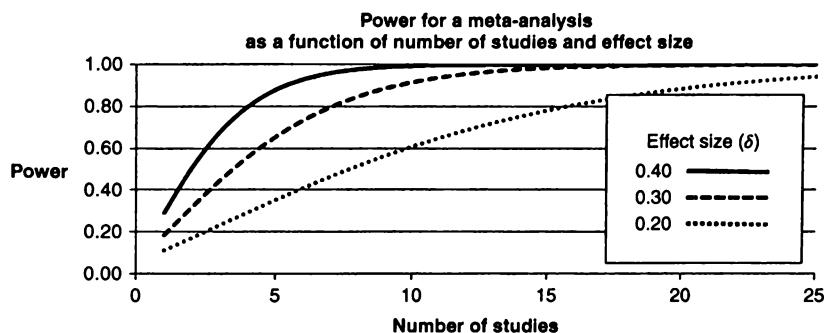


Figure 34.2 Power for a meta-analysis as a function of number of studies and  $\delta$ .

### Power for main effect using random-effects model

The formulas for significance and for power under the random-effects model have the same structure as those for a fixed effect meta-analysis. Using the notation introduced in Chapter 12 where an asterisk (\*) indicates that a statistic is based on random-effects variance, the test of significance for the main effect is still based on a test statistic (for example)  $Z^*$ , computed as

$$Z^* = \frac{M^*}{\sqrt{V_{M^*}}}, \quad (34.13)$$

but  $M^*$  and  $V_{M^*}$  are now the estimated mean effect size and its variance using random-effects weights. As before,  $Z^*$  is evaluated with reference to the standard normal distribution which yields the corresponding  $p$ -value. For a two-tailed test,

$$p^* = 2[1 - (\Phi(|Z^*|))]. \quad (34.14)$$

Similarly, power analysis is still based on lambda ( $\lambda^*$ ), defined as

$$\lambda^* = \frac{\delta^*}{\sqrt{V_{\delta^*}}}, \quad (34.15)$$

but now  $\delta^*$  and  $V_{\delta^*}$  are the true mean effect size and its variance for the summary effect. The variance incorporates variance within studies and variance between studies. Consider the simple situation in which each study has the same within-study variance, say  $V_y$ . Then the variance may be written as

$$V_{\delta^*} = \frac{V_y + \tau^2}{k}. \quad (34.16)$$

Plausible values of the within-study variance,  $V_y$ , might be obtained using the same procedures as those used for the fixed-effect model. Plausible values of the between-studies variance,  $\tau^2$ , might be obtained using data from the pilot study, by computing the effect sizes for the studies gathered as part of the pilot and looking at how much these effects actually vary from study to study. Alternatively, the between-studies variance in a previous, similar, meta-analysis might be suitable.

Finally, Hedges and Pigott propose a convention that can be used to represent small, medium, and large degrees of heterogeneity. This convention is to set  $\tau^2$  equal to 0.33, 0.67, or 1.0 times the within-study variance, so that the total variance  $V\delta^* = 1.33V_y/k$ ,  $1.67V_y/k$ , or  $2.00V_y/k$

As before, the critical value of alpha is given by

$$C_\alpha = \Phi(1 - \alpha/2), \quad (34.17)$$

and power for a two-tailed test is then given by

$$\text{Power} = 1 - \Phi(C_\alpha - \lambda^*) + \Phi(-C_\alpha - \lambda^*). \quad (34.18)$$

In practice, the effect size is likely to be the same under the random-effects model as it had been under the fixed-effect model. However, the variance will always be larger under the random-effects model as compared with the fixed-effect model (see Chapter 13).

### Illustrative example

The example above for the fixed-effect model can be used here as well if we assume that the true effect size varies from study to study, and the random-effects model is therefore appropriate. Again, we assume an effect size of  $d = 0.30$  and 10 studies with 25 patients per group. The within-study variance for one study is computed as

$$V_y = \frac{25 + 25}{25 \times 25} + \frac{0.30^2}{2(25 + 25)} = 0.0809.$$

If we assume a moderate degree of between-study heterogeneity and apply the convention, the variance of the summary effect is computed as

$$V_\delta^* = \frac{1.667 \times 0.0809}{10} = 0.0135.$$

The parameter  $\lambda^*$  is computed as

$$\lambda^* = \frac{0.30}{\sqrt{0.0135}} = 2.5836.$$

The critical value of alpha is given by

$$C_\alpha = \Phi(1 - 0.05/2) = 1.96.$$

In Excel, =NORMSINV(1-0.05/2) returns 1.96.

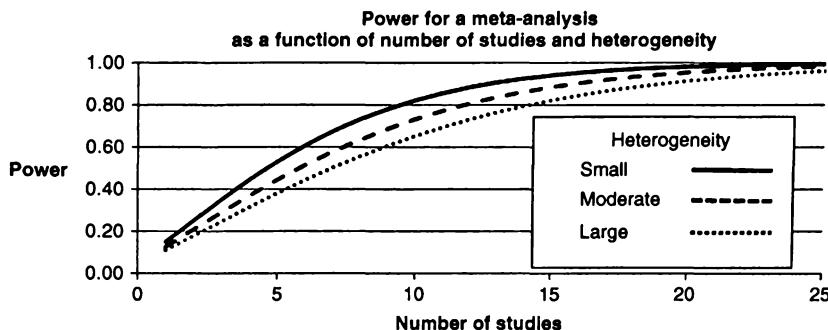
Power is then given by

$$\text{Power} = 1 - \Phi(1.96 - 2.5836) + \Phi(-1.96 - 2.5836) = 0.7336.$$

This computation can be done in EXCEL as

=1-NORMSDIST(1.96-2.5836) + NORMSDIST(-1.96-2.5836).

As before, we can create a graph that shows power under a series of assumptions. Figure 34.3 is based on the middle effect size of  $\delta = 0.30$ , and shows how power will vary if the dispersion is small, medium, or large. We could create a similar graph for an effect size ( $\delta$ ) of 0.20 or of 0.40.



**Figure 34.3** Power for a meta-analysis as a function of number of studies and heterogeneity.

Reading from left to right, if the between-studies dispersion is small, moderate, or large, then we would need about 12, 15, or 20 studies to get power of 90%. Assume that the pilot study shows that we can locate at least 15 studies. Reading from top to bottom at 15 on the  $x$ -axis, we see that if dispersion is small, medium or large, power would be around 0.95, 0.92 or 0.90. Of course, power will be lower for an effect size of 0.20 and higher for an effect size of 0.40.

### Notes about random effects

Under the random-effects model the variance of each study (and of the summary effect) includes two components, the variance within studies and the variance between studies. The practical impact is that power will depend both on the total sample size and also the number of studies. If there is substantial study-to-study dispersion, then the only way to yield good power is to include a large number of studies (which reduces this element of the variance). If the number of studies is small, then power could remain low even if the total sample size across studies reaches into the tens of thousands or even higher.

It is important to understand that the fixed-effect model and random-effects model address different hypotheses, and that they use different estimates of the variance because they make different assumptions about the nature of the distribution of effects across studies, as explained in Chapter 13. Researchers sometimes remark that power is lower under the random-effects model than for the fixed-effect model. While this statement may be true, it misses the larger point: it is not meaningful to compare power for fixed- and random-effects analyses since the two values of power are not addressing the same question.

### POWER ANALYSIS FOR A TEST OF HOMOGENEITY

Many meta-analyses include a test of homogeneity, which asks whether or not the between-studies dispersion is more than would be expected by chance. The test of

significance is discussed in Chapter 16, and is based on  $Q$ , the sum of the squared deviations of each study's effect size estimate ( $Y_i$ ) from the summary effect ( $M$ ), with each deviation weighted by the inverse of that study's variance. Concretely,

$$Q = \sum_{i=1}^k W_i(Y_i - M)^2. \quad (34.19)$$

$Q$  is then evaluated with reference to the chi-squared distribution with  $k - 1$  degrees of freedom. Power for this test depends on three factors. The larger the ratio of between-studies to within-studies variance, the larger the number of studies, and the more liberal the criterion for significance, the higher the power. We can compute power under two alternative scenarios. In the first scenario we do not assume any meta-analytic model (such as fixed-effect or random-effects) for the true effects in the different studies. In the second scenario we assume a random-effects model for these true effects. In both scenarios we assume that the within-study variance  $V_{Y_i}$  is the same in each study. For a discussion of power analysis when the  $V_{Y_i}$  are not equal, see Hedges and Pigott (2001).

### Power for a test of homogeneity in the absence of a meta-analytic model

Researchers who apply the fixed-effect model sometimes test the dispersion for significance. Technically, if the fixed-effect model of homogeneity is true, there is no dispersion, since the model asserts that the between-studies variance is zero. However, it is possible to perform this test (and to compute power) in the absence of a meta-analytic model (a model in which the effect sizes are taken to be fixed, but not necessarily equal).

The expected value of the test statistic ( $Q$ ) is equal to  $df + \lambda$ , where the noncentrality parameter  $\lambda$  is computed as

$$\lambda = df \times \left( \frac{\tau^2}{V_Y} \right). \quad (34.20)$$

In this formula  $df$  is  $k-1$  ( $k$  being the number of studies). Here  $\tau^2$  and  $V_Y$  are the between-studies variance and the within-studies variance respectively, but power depends only the ratio of these two values (rather than their absolute values).

Power is then given by

$$\text{Power} = 1 - F(C_\alpha | k - 1; \lambda), \quad (34.21)$$

where

$$F(X | df; \lambda)$$

is the cumulative distribution function of a noncentral chi-square with  $df$  degrees of freedom and noncentrality parameter  $\lambda$ , and where  $C_\alpha$  is the  $100(1 - \alpha)$  percent point of the central chi-squared distribution.

To approximate this value in Excel we can proceed as follows. First, the function =CHIINV(*alpha*, *df*) returns the alpha critical value  $C_\alpha$ . Then,

$$a = 1 + \frac{\lambda}{df + \lambda}, \quad (34.22)$$

$$b = df + \frac{\lambda^2}{df + 2\lambda}, \text{ and} \quad (34.23)$$

$$X = \frac{C_\alpha}{a}. \quad (34.24)$$

Finally, the expression =1-GAMMADIST(*X*, *b*/2, 2, TRUE) returns the value of power.

### **Illustrative example**

Suppose that we are planning a meta-analysis with six studies, and we want to compute power to test for a large amount of dispersion (in which  $\tau^2 = V_Y$ ) with alpha set at 0.05. Then,

$$\lambda = 5 \left( \frac{1}{1} \right) = 5.$$

(Note that we don't need to know the actual variances, only the ratio of  $\tau^2$  to  $V_Y$ . Using the conventions, this ratio is 1.000., 0.667, or 0.333 for large, moderate, and small dispersion.)

In Excel, the chi-squared critical value corresponding to  $\alpha$  of 0.05 for 5 *df* is given by the function =CHIINV(0.05,5), which returns 11.0705. The intermediate values *a* and *b* are computed as

$$a = 1 + \frac{5}{5+5} = 1.5 \quad \text{and}$$

$$b = 5 + \frac{5^2}{5+2\times 5} = 6.6667.$$

Finally, power is given by the Excel function =1-GAMMADIST(11.0705/1.5, 6.6667/2, 2, TRUE), which returns 0.3553.

Note that we assumed a large amount of dispersion, and the number of studies is typical for many meta-analyses, yet power is quite low. In this example we would need 26 studies to boost power into the 80% range.

### **Power for a test of homogeneity under the random-effects model**

Power of the test is slightly different if we impose a random-effects model on the effects across the studies. The formulas are the same as those for the test in the absence of a meta-analytic model except for the final step. Now, we use instead an Excel expression of the form =CHIDIST(*x*, *df*) to return power, where

$$X = \frac{C_\alpha}{1 + \frac{\tau^2}{V_Y}}. \quad (34.25)$$

In this example, =CHIDIST(11.0705/(1 + 1/1); 5) returns 0.3541, or 35.4%.

In this chapter we have presented the logic of power analysis for systematic reviews, and presented formulas for tests of the main effect and for tests of homogeneity. Hedges and Pigott (2004) discuss power for tests of subgroup analyses, and for meta-regression.

### SUMMARY POINTS

- The process of power analysis for a meta-analysis closely parallels the issues of power analysis for a primary study.
- Under the fixed-effect model, the 'sample size' factor is driven by the number of subjects accumulated across studies. Power for a meta-analysis of ten studies with 100 persons each is the same as power for one study with 1000 persons (provided the effect size is constant). Therefore, a fixed-effect analyses with a decent number of studies and/or some large studies will often have good power to detect any nontrivial effect size.
- Under the random-effects model power depends not only on the total number of subjects but also on the number of studies. Even if the effect size is large and the cumulative number of subjects is large, if there are only a few studies and the variance between-studies is substantial, power could be very low.
- The absence of statistical significance should never be interpreted as evidence that an effect is absent. This is important to keep in mind since power to detect heterogeneity in effect sizes, and power to detect the relationship between subgroup membership and effect size, or between covariate values and effect size, is often quite low.

### Further Reading

- Borenstein, M. (1994). The case for confidence intervals in controlled clinical trials. *Controlled Clinical Trials, Oct*; 15(5): 411–428.
- Borenstein, M. (1997). Hypothesis testing and effect size estimation in clinical trials. *Annals of Allergy, Asthma and Immunology, Jan*; 78(1): 5–16(12).
- Cohen, J. (1987). *Statistical Power Analysis for the Behavioral Sciences* (2<sup>nd</sup> edn). Hillsdale, NJ: Lawrence Erlbaum Associates.
- Hedges L.V. and Pigott T.D. (2001). The power of statistical tests in meta-analysis, *Psychological Methods* 6: 203–217.
- Hedges L.V. and Pigott T.D. (2004). The power of statistical tests for moderators in meta-analysis, *Psychological Methods* 9: 426–445.
- Sutton, A.J., Cooper, N.J., Jones, D.R., Abrams, K.A., Lambert, P., Thompson, J.R. (2007). Evidence based sample size calculations for the future trials based on results of current meta-analysis. *Statistics in Medicine* 26: 2479–2500.
- Sutton, A.J., Donegan, S., Takwoingi, Y., Garner, P., Gamble, C., Donald, A. (2009) An encouraging assessment of methods to inform priorities for updating systematic reviews. *Journal of Clinical Epidemiology, Mar*; 62(3): 241–251.

# Publication Bias

---

### Introduction

The problem of missing studies

Methods for addressing bias

Illustrative example

The model

Getting a sense of the data

Is there evidence of any bias?

How much of an impact might the bias have?

Summary of the findings for the illustrative example

Conflating bias with the small-study effect

USING logic to disentangle bias from small-study effects

These methods do not give us the 'correct' effect size

Some important caveats

Procedures do not apply to studies of prevalence

The model for publication bias is simplistic

Concluding remarks

Putting it all together

---

## INTRODUCTION

While a meta-analysis will yield a mathematically accurate synthesis of the studies included in the analysis, if these studies are a biased sample of all relevant studies that had been performed, the mean effect computed by the meta-analysis will reflect this bias. Several lines of evidence show that studies that report relatively high effect sizes are more likely to be published than studies that report lower effect sizes. Since published studies are more likely to find their way into a meta-analysis, any bias in the literature is likely to be reflected in the meta-analysis as well. This issue is generally known as publication bias.

The problem of publication bias is not unique to systematic reviews. It affects the researcher who writes a narrative review and even the clinician who is searching a database for primary papers. Nevertheless, it has received more attention with regard

to systematic reviews and meta-analyses, possibly because these are promoted as being more accurate than other approaches to synthesizing research.

This chapter is structured as follows:

- We discuss the reasons for publication bias and the evidence that it exists. We also outline a series of methods that have been developed to assess the likely impact of bias in any given meta-analysis.
- We introduce the idea of a small-study effect, and how this is often conflated with publication bias. In particular, we explain that the methods for assessing publication bias should be seen as a sensitivity analysis, rather than an attempt to compute the ‘correct’ effect size.
- We discuss the conditions that must be met before we can apply any of the methods addressed in this chapter.

## THE PROBLEM OF MISSING STUDIES

When planning a systematic review, we develop a set of inclusion criteria that govern the types of studies that we want to include. Ideally, we would be able to locate all studies that meet our criteria, but in the real world, this is rarely possible. Even with the advent of (and perhaps partly due to an over-reliance on) electronic searching, it is likely that some studies which meet our criteria will escape our search and not be included in the analysis.

If the missing studies are a *random* subset of all relevant studies, the failure to include these studies will result in less information, wider confidence intervals, and less powerful tests, but would have no systematic impact on the effect size. However, if the missing studies are *systematically* different than the ones we were able to locate, then our sample will be biased. The specific concern is that studies that report relatively large effects for a given question are more likely to be published than studies that report smaller effects for the same question. This leads to a bias in the published literature, which then carries over to a meta-analysis that draws on this literature.

### Studies with significant results are more likely to be published

Several lines of research (reviewed by Dickersin, 2005) have established that studies with statistically significant results are more likely to find their way into the published literature than studies that report results that are not statistically significant. And, for any given sample size, the result is more likely to be statistically significant if the effect size is larger. It follows that if there is a universe of studies that looked at the magnitude of a relationship, and the observed effects are distributed over a range of values (as they always are), the studies with effects toward the higher end of that range are more likely to be statistically significant and therefore to be published. This tendency has the potential to produce very large biases in the magnitude of the relationships, particularly if studies have relatively small sample sizes (see Hedges, 1984, 1989).

A particularly enlightening line of research was to identify groups of studies as they were initiated and then follow them prospectively over a period of years to see which were published and which were not. This approach was taken by Easterbrook, Berlin, Gopalan, and Matthews (1991), Dickersin, Min, and Meinert (1992), and Dickersin and Min (1993a), among others. Nonsignificant studies were less likely to be published than significant studies (61–86% as likely) and when published were subjected to longer delay prior to publication. Similar investigations have demonstrated that researchers selectively report their findings in the reports they do publish, sometimes even changing what is labeled *a priori* as the main hypothesis (Chan *et al.*, 2004).

### **Published studies are more likely to be included in a meta-analysis**

If persons performing a systematic review were able to locate studies that had been published in the grey literature (any literature produced in electronic or print format that is not controlled by commercial publishers, such as technical reports and similar sources), then the fact that the studies with higher effects are more likely to be published in the more mainstream publications would not be a problem for meta-analysis. In fact, though, this is not usually the case.

While a systematic review *should* include a thorough search for all relevant studies, the actual amount of grey/unpublished literature included, and the types vary considerably across meta-analyses. When Rothstein (2006) reviewed the 95 meta-analytic reviews published in *Psychological Bulletin* between 1995 and 2005 to see whether they included unpublished or grey research, she found that 23 of the 95 clearly did not include any unpublished data. Clarke and Clarke (2000) studied the references from healthcare protocols and reviews published in The Cochrane Library in 1999 and found that about 92% of references to studies included in reviews were to journal articles. Of the remaining 8%, about 4% were to conference proceedings, about 2% were to unpublished material (for example personal communication, *in press* documents and data on file), and slightly over 1% were to books or book chapters. In a similar vein, Mallet, Hopewell, and Clarke (2002) looked at the sources of grey literature included in the first 1000 Cochrane systematic reviews and found that nearly half of them did not include any data from grey or unpublished sources. Since the meta-analyses published in the Cochrane Database have been shown to retrieve a higher proportion of studies than those published in many journals, these estimates probably underestimate the extent of the problem.

Some have suggested that it is legitimate to exclude studies that have not been published in peer-reviewed journals because these studies tend to be of lower quality. For example, in their systematic review, Weisz *et al.* (1995) wrote ‘We included only published psychotherapy outcome studies, relying on the journal review process as one step of quality control’ (p. 452). However, it is not obvious that journal review assures high quality, nor that it is the *only* mechanism that can do so. For one thing, not all researchers aim to publish their research in academic journals. For example, researchers working for government agencies, independent think-tanks, or consulting firms generally focus on producing reports, not journal articles. Similarly, a thesis or

dissertation may be of high quality, but is unlikely to be submitted for publication in an academic journal if the individual who produced it is not pursuing an academic career. And of course, peer review may be biased, unreliable, or of uneven quality. Overall, then, publication status cannot be used as a proxy for quality and in our opinion should not be used as a basis for inclusion or exclusion of studies.

Some researchers think that publication bias refers to a distinction between studies that were published in a journal vs. those that were published in the grey literature (as technical reports, dissertations, or abstracts). This is incorrect. Publication bias (or more generally, retrieval bias) refers to fact that we can locate some studies and not others. In this context, all studies that we can locate are considered published (or retrieved), without regard to where they were located.

### Other sources of bias

Other factors that can lead to an upward bias in effect size and are included under the umbrella of publication bias are the following. Language bias (English-language databases and journals are more likely to be searched, which leads to an oversampling of statistically significant studies) (Egger *et al.*, 1997; Jüni *et al.*, 2002); availability bias (selective inclusion of studies that are easily accessible to the researcher); cost bias (selective inclusion of studies that are available free or at low cost); familiarity bias (selective inclusion of studies only from one's own discipline); duplication bias (studies with statistically significant results are more likely to be published more than once [Tramer *et al.*, 1997]); and citation bias (whereby studies with statistically significant results are more likely to be cited by others and therefore easier to identify [Gøtzsche, 1987; Ravnklov, 1992]).

## METHODS FOR ADDRESSING BIAS

In sum, it is possible that the studies in a meta-analysis may overestimate the true effect size because they are based on a biased sample of the target population of studies. But how do we deal with this concern?

The best approach would be for the reviewer to perform a truly comprehensive search for studies, in hopes of minimizing the bias. Despite the increased resources that are needed to locate and retrieve data from sources such as dissertations, theses, conference papers, government and technical reports, and the like, it is generally indefensible to conduct a synthesis that categorically excludes these types of research reports. Potential benefits and costs of grey literature searches must be balanced against each other. Readers who would like more guidance in the process of literature searching and information retrieval may wish to consult Hopewell, Mallett and Clarke (2005), Reed and Baxter (2009), Rothstein and Hopewell (2009), Wade, Turner, Rothstein and Lavenberg (2006), Giustinti (2019), Glanville (2019), Lefebvre *et al.* (2019).

There is evidence that this approach is somewhat effective. Cochrane reviews tend to include more studies and to report a smaller effect size than similar reviews published in medical journals. Serious efforts to find unpublished, and *difficult to find*

studies, typical of Cochrane reviews, may therefore reduce some of the effects of publication bias.

While this approach may mitigate the impact of bias, it is likely that some bias will remain. Researchers have developed methods intended to assess its potential impact on any given meta-analysis. We shall illustrate these methods as they apply to a meta-analysis on the relationship between passive smoking and lung cancer.

### ILLUSTRATIVE EXAMPLE

Hackshaw *et al.* (1997) published a meta-analysis with data from 37 studies that reported on the relationship between so-called second-hand (passive) smoking and lung cancer. The paper reported that exposure to second-hand smoke increased the risk of lung cancer in the nonsmoking spouse by about 24%. Questions were raised about the possibility that studies with larger effects were more likely to have been published (and included in the analysis) than those with smaller (or nil) effects and that the conclusion was therefore suspect.

### THE MODEL

In order to gauge the impact of publication bias, we need a model that tells us which studies are likely to be missing. The model that is generally used (and the one we follow here) makes the following assumptions: (a) Large studies are likely to be published regardless of statistical significance because these involve large commitments of time and resources; (b) moderately sized studies are at risk for being lost, but with a moderate sample size even modest effects will be significant, and so only some studies are lost here; and (c) small studies are at greatest risk for being lost. Because of the small sample size, only the larger effects are likely to be significant, with the smaller effects less likely to be published.

The combined result of these three items is that we expect the bias to increase as the sample size goes down, and the methods described below are all based on this model. Other, more sophisticated methods have been developed for estimating the number of missing studies and/or adjusting the observed effect to account for bias. These have rarely been used in actual research because they are difficult to implement and also because they require the user to make some relatively sophisticated assumptions and choices.

Before proceeding with the example, we call the reader's attention to an important caveat. The procedures that we describe here look for a relationship between sample size and effect size, and if such a relationship is found, it is attributed to the existence of missing studies and is assumed to introduce bias. In the interest of readability, we will use the term *bias* when discussing these methods. However, while publication bias is *one possible reason* that the effect size is larger in the smaller studies, it is also possible that the effect size *really is larger* in the smaller studies. We return to this later in this chapter, in the section entitled *Conflating bias with the small-study effect*.

## GETTING A SENSE OF THE DATA

A good place to start in assessing the potential for bias is to get a sense of the data, and the forest plot can be used for this purpose. Figure 35.1 is a forest plot of the studies in the passive smoking meta-analysis. In this example, an increase in risk is indicated by a risk ratio greater than 1.0. The overwhelming majority of studies show an increased risk for second-hand smoke, and the last row in the spreadsheet shows the summary data for the random-effects model. The risk ratio is 1.238, and the 95% confidence interval is 1.129 to 1.356.

The studies have been plotted from most precise to least precise, so that larger studies appear toward the top and smaller studies appear toward the bottom. This has no impact on the summary effect, but it allows us to see the relationship between sample size and effect size. There does seem to be a tendency for smaller studies to have larger effects. For example, among the larger studies (the first eighteen on the plot) there are none with risk ratios higher than 2.0. By contrast, among the nineteen smaller studies there are eight studies with risk ratios higher than 2.0. This is not intended as a statistical index, but simply to provide some context for what follows.

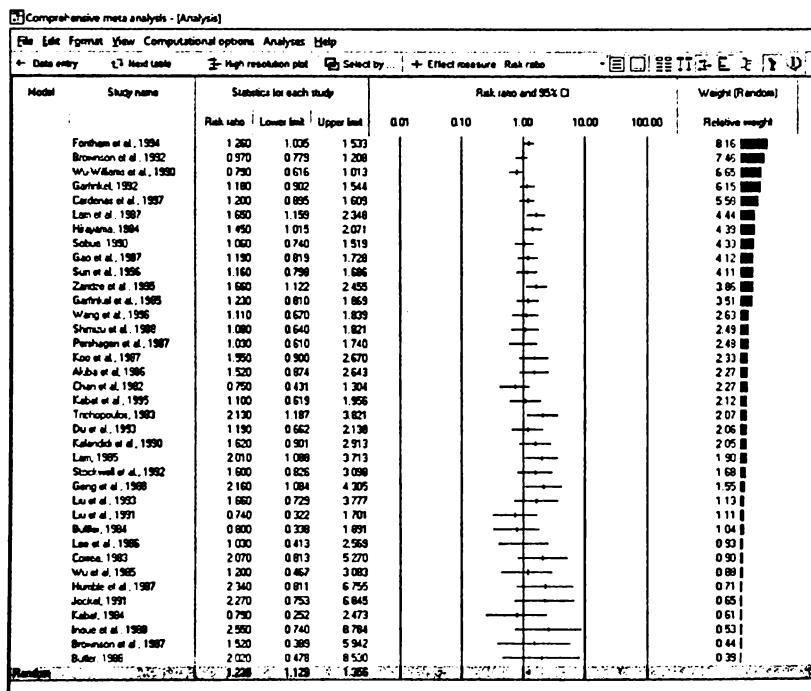
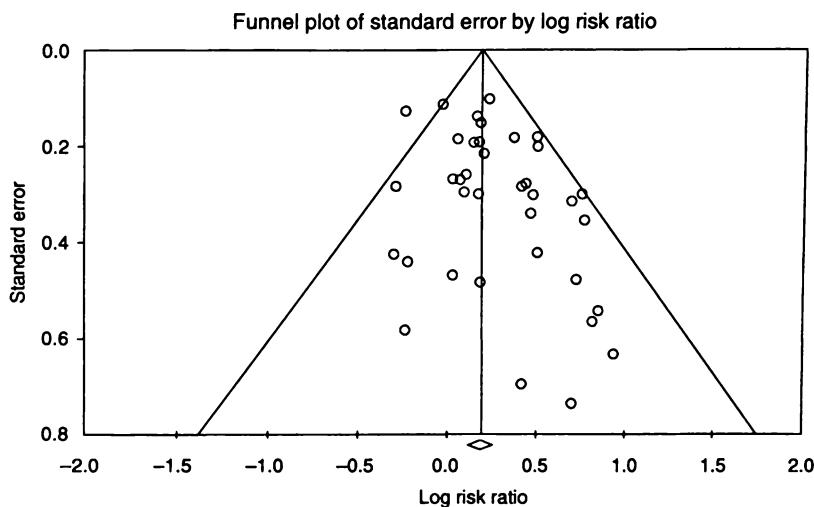


Figure 35.1 Passive smoking and lung cancer – forest plot.



**Figure 35.2** Passive smoking and lung cancer – funnel plot.

### The funnel plot

Another mechanism for displaying the relationship between study size and effect size is the funnel plot. Figure 35.2 shows the funnel plot for the smoking analysis.

In this figure, the effect size is plotted on the  $X$ -axis. Since we are working with risk ratios, we plot this in log units. A log value of 0.0 corresponds to a risk ratio of 1.0. The mean value of 0.213 (in log units) corresponds to the risk ratio of 1.238 (in risk ratio units) in Figure 35.1. The  $y$ -axis is the standard error of the effect size in each study. The scale is reversed, with low standard error (large studies) at the top and high standard error (small studies) at the bottom.

In the absence of bias, we would expect roughly half the studies to fall on either side of the mean. On the other hand, the model for bias we outlined a moment ago predicts a different pattern. Toward the top of the plot, where the studies are large, we would expect to see roughly the same number of studies on either side of the mean since most studies are published. By contrast, toward the lower parts of the plot, where the studies are smaller, we would expect to see more studies toward the right and fewer toward the left, since those on the left are less likely to be statistically significant and therefore more likely to be missing.

The pattern in Figure 35.2 does seem to resemble the pattern we associate with bias. Toward the bottom of the plot, there are a number of studies on the right, but relatively few studies on the left. It is possible that the studies on the left were performed, but are missing from our analysis.

Tangentially, this is referred to as a funnel plot since an earlier version of this plot employed a different index for the  $y$ -axis and resembled a funnel (Light and Pillemer, 1984; Light *et al.*, 1994).

## IS THERE EVIDENCE OF ANY BIAS?

We suggested above that there appears to be a relationship between study size and effect size. Because the interpretation of a funnel plot is largely subjective and unreliable (Terrin, Schmid, and Lau, 2005), several tests have been proposed to quantify or test the relationship between sample size and effect size. Among these, Begg and Mazumdar (1994) proposed a rank correlation test. This computes the rank correlation between the effect size and the standard error. In Figure 35.1 and in Figure 35.2, if we generated a line based on this correlation, this line would extend from the top left toward the bottom right. Similarly, Egger, Davey Smith, Schneider, and Minder (1997) proposed a regression test. This takes the same approach as Begg and Mazumdar but uses regression rather than a rank correlation and has better statistical power. These methods are discussed in Higgins and Thomas (2019), Rothstein, Sutton, and Borenstein (2005), and Vevea, Coburn, and Sutton (2019).

In the passive smoking analysis, the rank correlation is 0.14 and the one-tailed  $p$ -value is 0.107. The Egger regression yields an intercept of 0.89 and a one-tailed  $p$ -value of 0.012.

## HOW MUCH OF AN IMPACT MIGHT THE BIAS HAVE?

While the tests outlined above are widely reported, they have limited utility. Even if these tests work properly and if a statistically significant result really is indicative of bias, the tests only tell us that bias exists. This is not terribly informative, since we can assume that bias exists based on the research outlined earlier. The more relevant question is how much of an impact the bias might have had on our analysis, and these methods do not directly address this question.

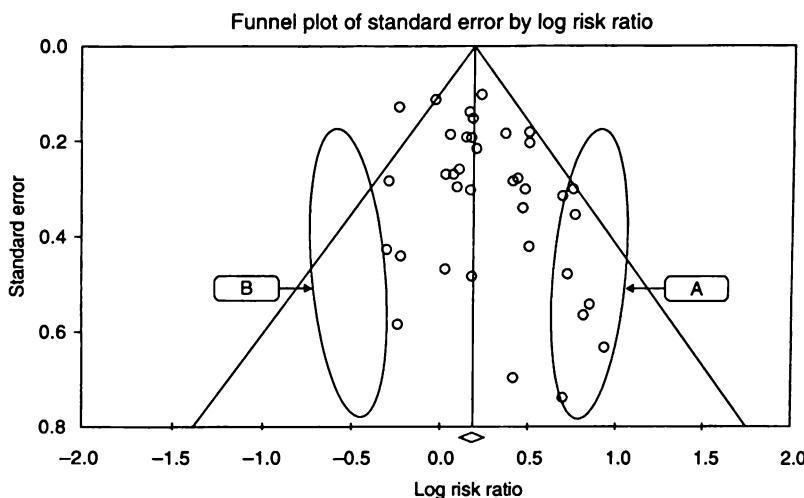
A more informative approach would be to quantify the potential impact of the bias. Based on this, we might to classify each meta-analysis into one of three broad groups, as follows.

- The impact of bias is probably trivial. If all relevant studies were included, the effect size would probably remain largely unchanged.
- The impact of bias is probably modest. If all relevant studies were included, the effect size might shift but the key finding (that the effect is, or is not, of substantive importance) would probably remain unchanged.
- The impact of bias may be substantial. If all relevant studies were included, the key finding (that the effect size is, or is not, of substantive importance) could change.

Some procedures that may allow us to classify the results in this way are described immediately below.

### Duval and Tweedie's *Trim and Fill*

The key idea behind the funnel plot is that publication bias may be expected to lead to asymmetry. Toward the bottom of Figure 35.2, there are studies on the right but not the left. We are concerned that the studies on the left were actually performed,



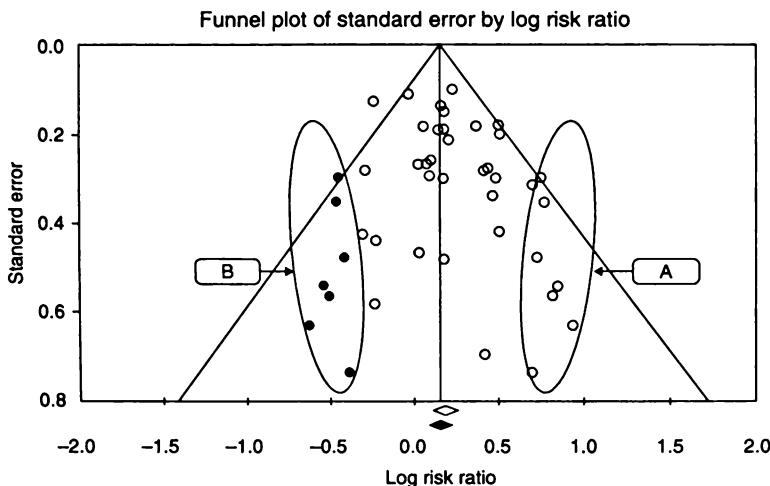
**Figure 35.3** Observed studies only.

but are missing from the analysis. The approach taken by Trim and Fill is to estimate what studies might be missing and then add them to the analysis (Duval and Tweedie, 2000a, 2000b).

Figure 35.3 is the same as Figure 35.2, but includes two annotations. The ellipse [A] denotes the studies Trim and Fill identifies as contributing to the asymmetry. The ellipse [B] corresponds to the gap on the other side. In the absence of bias, we would expect to see studies in this section. Figure 35.4 shows the results of the Trim and Fill analysis. The algorithm has created the missing studies and inserted them into the analysis. The new studies are shown as closed circles.

At the bottom of Figure 35.4, we see the results of the analysis based on the observed studies (open diamond) as well as the observed plus imputed studies (filled diamond). The observed point estimate in log units is 0.213, corresponding to a risk ratio of 1.238. The imputed point estimate in log units is 0.173, corresponding to a risk ratio of 1.189.

The imputed point estimate falls to the left of the original point estimate. This suggests that the original estimate was too high. For our purposes, the more important point is that the imputed point estimate is fairly close to the original. The original risk ratio was 1.238, and the adjusted risk ratio is 1.189. In this context, these two estimates have substantially the same meaning. That is, someone who is concerned that their risk of developing lung cancer is increased by 24% will also be concerned if their risk is increased by 19%. Using the classification system we proposed earlier we would conclude that bias may have led us to overestimate the risk, but the substantive conclusions of the analysis are valid. By contrast, if the imputed value suggested that smoking increased the risk by 2% (for example), we would say that the results were not robust.



**Figure 35.4** Observed studies and studies imputed by Trim and Fill.

The Trim and Fill procedure is built entirely on the assumption that asymmetry must be due to publication bias. In fact, though, asymmetry may be due to other causes entirely. This is an issue that affects all the procedures outlined in this chapter, and so we will address it in a separate section toward the end of this chapter.

### Restricting analysis to the larger studies

If publication bias operates primarily on smaller studies, then restricting the analysis to the larger studies, which might be expected to be published irrespective of their results, would in theory mitigate the extent of the problem. Rather than draw a dichotomy between large and small studies, we can work with study size as a continuous factor and perform a cumulative meta-analysis.

A cumulative meta-analysis is a meta-analysis run first with one study and then repeated with a second study added, then a third, and so on. Similarly, in a cumulative forest plot, the first row shows the effect based on one study, the second row shows the cumulative effect based on two studies, and so on.

Figure 35.5 shows a cumulative forest plot of the data. This is based on the same studies as Figure 35.1, and (as before) the studies have been sorted from the most precise to the least precise (roughly corresponding to largest to smallest). However, this plot reflects a cumulative analysis. Here, the first row is a 'meta'-analysis based only on the Fontham *et al.* study. The second row is a meta-analysis based on two studies (Fontham *et al.* and Brownson *et al.*), and so on. The last study to be added is Butler (1988), and so the point estimate and confidence interval shown on the line labeled 'Butler' are identical to that shown for the summary effect. The scale on

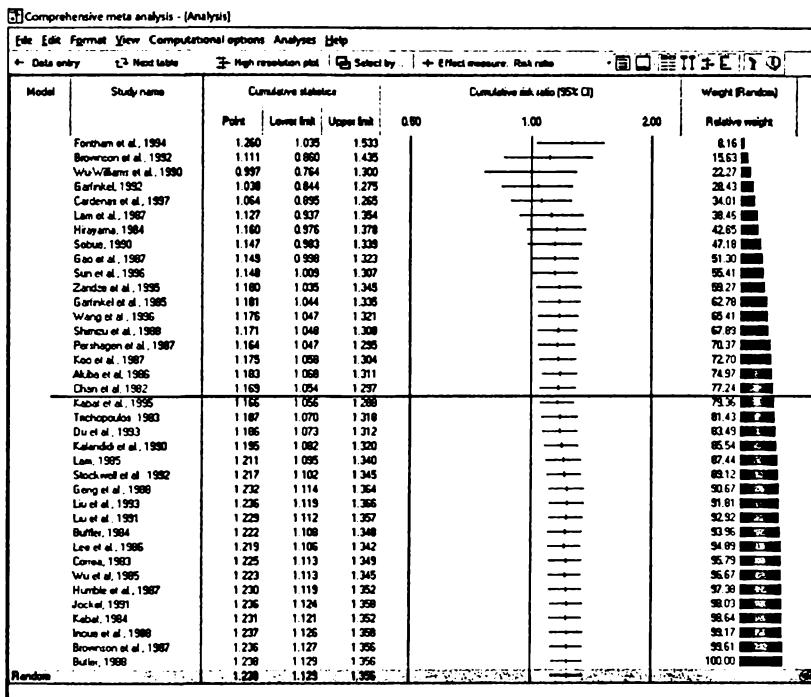


Figure 35.5 Passive smoking and lung cancer – cumulative forest plot.

this plot is 0.50 to 2.00 rather than 0.01 to 100 as in the prior plot. For purposes of this discussion, we divide the studies into the 18 larger ones and the 19 smaller ones, as indicated by a horizontal line. The selection of this specific dividing point is arbitrary.

If the point estimate has stabilized with the inclusion of the larger studies and does not shift with the addition of smaller studies, then there is no reason to assume that the inclusion of smaller studies had injected a bias (i.e. since it is the smaller studies in which study selection is likely to be greatest). On the other hand, if the point estimate does shift when the smaller studies are added, then there is at least a *prima facie* case for bias, and one would want to investigate the reason for the shift.

With the 18 largest studies in the analysis, starting at the top (inclusive of Chan & Fung, 1981) the cumulative relative risk is 1.169. With the addition of the remaining 19 smaller studies, the point estimate shifts to the right, and the relative risk is 1.238. As such, our estimate of the relative risk has increased. However, the key point is that even if we had limited the analysis to the 18 larger studies, the relative risk would have been 1.169 (with 95% confidence interval of 1.054, 1.297) and the clinical implications probably would have been the same.

## Egger's regression

A third option for estimating the *potential* impact of bias is Egger's regression. While this test is typically used to determine whether there is evidence of any bias, the intercept in this regression can also be used to estimate the effect size that we would see in a large study. Variants of this approach are discussed in Moreno, Sutton, Ades, Cooper, and Abrams (2011), Moreno, Sutton, Ades, *et al.* (2009), Moreno *et al.* (2012), Moreno, Sutton, Turner, *et al.* (2009), Page, Higgins, and Sterne (2019), Rucker, Carpenter, and Schwarzer (2011), and Rucker, Schwarzer, Carpenter, Binder, and Schumacher (2011).

## SUMMARY OF THE FINDINGS FOR THE ILLUSTRATIVE EXAMPLE

### Getting a sense of the data

There are a substantial number of studies in this analysis. While the vast majority show an increased risk, only a few are statistically significant. This suggests that the mechanism of publication bias based on statistical significance was not a powerful one in this case.

### Is there evidence of bias?

The funnel plot is noticeably asymmetric, with a majority of the smaller studies clustering to the right of the mean. This visual impression is confirmed by Egger's test which yields a statistically significant *p*-value. The rank correlation test did not yield a significant *p*-value, but this could be due to the low power of the test. As a whole, the smaller studies did tend to report a higher association between passive smoking and lung cancer than did the larger studies.

### What impact might the bias have on the risk ratio?

The complete meta-analysis showed that passive smoking was associated with a 24% increase in risk of lung cancer. The Trim and Fill algorithm tells us that the actual risk increase might be 19%. The cumulative analysis tells us that the actual risk increase might be closer to 17%. Critically, the substantive implications for all of these values are similar. Earlier, we suggested that the goal of a publication bias analysis should be to classify the results into one of three categories (a) where the impact of bias is trivial, (b) where the impact is not trivial but the major finding is still valid, and (c) where the major finding might be called into question. This meta-analysis seems to fall squarely within category *b*. There *is* evidence of larger effects in the smaller studies, which is consistent with our model for publication bias. However, there is no reason to doubt the validity of the core finding, that passive smoking is associated with a clinically important increase in the risk of lung cancer.

## CONFLATING BIAS WITH THE SMALL-STUDY EFFECT

To this point, we have outlined two methods that test for the presence of bias and two methods that estimate the potential impact of the bias. In fact, none of these methods directly assess bias. Rather, the methods all look at the relationship between effect size and study size. If they find this relationship, they attribute it to publication bias. However, this attribution may not be correct. If the effect size is larger in the smaller studies, this could be due to bias, but could also be due to other factors entirely.

The idea that we can identify publication bias by looking for a larger effect size in smaller studies works well in cases where the fixed-effect model is called for – that is, when all studies are estimating a common parameter. In that case, when the effect size is larger in smaller studies, there are only two possible reasons:

- A. Random sampling error
- B. Publication bias

If the test for asymmetry is statistically significant, we can rule out (A) and conclude that publication bias exists (B). If we are using Trim and Fill or the cumulative procedure, we can estimate the impact of the bias.

However, when we are locating studies in the literature, we are rarely dealing with a case where all studies are estimating a common parameter. Rather, we assume that the true effect size varies from study to study. In this case, if the effect size is larger in smaller studies, there are three possible reasons:

- A. Random sampling error
- B. Publication bias
- C. The effect size really is larger in the smaller studies

If the test for asymmetry is statistically significant, we can rule out (A), but we cannot distinguish between (B) and (C). Similarly, the Trim and Fill procedure and cumulative analysis can tell us that the effect size is larger in smaller studies but cannot tell us if this is because of (B) or (C).

In fact, there are any number of reasons why the effect size could actually be larger in smaller studies. Consider the following examples.

1. Suppose that a new intervention is being studied. The initial trials are small, enroll patients who are very ill, and show large benefits from the treatment. Later trials are larger, enroll patients who are only moderately ill, and show more modest benefits. The effect size actually is larger in the smaller studies because the patients in these studies have more room to improve than those in the larger studies (Glasziou & Irwig, 1995; Smith & Egger, 1994; Stuck, Siu, Wieland, Adams, & Rubenstein, 1993).
2. Suppose that a new intervention is being studied. The initial trials are small and run by people who ensure that the patients take the medication as prescribed. Later trials are large and run by staff who are not able to track the patients as carefully. The effect size is larger in the smaller studies because the treatment in these studies is applied more consistently (Stuck, Rubenstein, & Wieland, 1998).

3. Suppose that an intervention is tested in a series of studies which vary in size. Large studies tend to be run by professionals who employ methods to minimize the risk of bias. Smaller studies are run by researchers with less experience and methodological flaws in these studies (for example, patients with a better prognosis being pushed into the treatment group) yield larger effects (Egger, Juni, Bartlett, Holenstein, & Sterne, 2003; Ioannidis, 2008b; Linde *et al.*, 1999; Terrin, Schmid, Lau, & Olkin, 2003; Wood *et al.*, 2008).
4. Suppose that researchers who plan the studies use a statistical power analysis to determine the sample size for each study. Suppose further that some studies are expected to yield a large effect (based on the population or the dose) and others are expected to yield a smaller effect. Based on the power analysis, the studies where we expect to see a large effect size will be assigned a smaller sample size and the studies where we expect to see a small effect size will be assigned a larger samples size. This would create the same pattern of results we expect to see based on the publication bias model (Linde *et al.*, 1997; Terrin *et al.*, 2003).

Sterne *et al.* (2001a) use the term *small-study effect* to describe a pattern where the effect is larger in small studies and to highlight the fact that the mechanism for this effect is not known. Using this terminology, we would say that there is clear evidence for a small-study effect. The fact that the effect size is larger in the smaller studies could reflect the presence of publication bias, but could also reflect the fact that the effect size actually is larger in smaller studies. Of course, both of these could contribute to the small-study effect (Ioannidis & Trikalinos, 2007; Peters *et al.*, 2010; Sterne, Egger, & Moher, 2008; Sterne *et al.*, 2011; Sterne & Egger, 2001; Sterne, Gavaghan, & Egger, 2000).

Critically, the small-study effect is simply one expression of heterogeneity. The effect size is smaller in some studies and larger in others, based on a myriad of factors. If the true effect size is larger in small studies, that is because of some factors that tend to be associated with these small studies, just as other factors may be associated with larger studies. If we are working with a random-effects analysis, all of these studies are part of the universe to which we are making an inference. The small studies cannot be said to inappropriately bias the mean effect upward any more than the large studies can be said to inappropriately bias the mean effect downward.

## USING LOGIC TO DISENTANGLE BIAS FROM SMALL-STUDY EFFECTS

Once we have ruled out random sampling error, we might be able to argue that the small-study effect is (or is not) probably due to publication bias (Sterne *et al.*, 2008; Sterne *et al.*, 2011).

For example, if most studies in the analysis are statistically significant and the effect size is larger in the smaller studies, publication bias is a plausible explanation. By contrast, if only a small proportion of studies are statistically significant, it is less likely that publication bias had a substantial impact on which studies were included.

The way studies were located might also be relevant in this context. If we are pulling studies from the literature that assessed the impact of drugs for treating depression, it is plausible to expect some bias based on selective publication and reporting. By contrast, in a prospective meta-analysis (where a set of primary studies had been planned in advance by a group of researchers and are now being included in a meta-analysis), we know that we have included all the trials, and so a small-study effect cannot be due to publication bias.

Another tool that may be used here is the contour-enhanced funnel plot (Peters, Sutton, Jones, Abrams, & Rushton, 2008). This is similar to the standard funnel plot, but allows us to distinguish between studies where the effect is small as opposed to studies where the effect is not statistically significant. For a discussion, see Page *et al.* (2019).

These examples are not intended to be exhaustive, but to provide a framework for thinking about possible causes of asymmetry in the funnel plot.

### THESE METHODS DO NOT GIVE US THE ‘CORRECT’ EFFECT SIZE

Since there is no way to distinguish between a small-study effect and publication bias, the methods which yield an adjusted effect size must be seen as a sensitivity analysis rather than seen as estimating the correct effect size. We can use these methods to say that:

- The impact of bias is probably trivial. If all relevant studies were included, the effect size would probably remain largely unchanged.
- The impact of bias is probably modest. If all relevant studies were included, the effect size might shift but the key finding (that the effect is, or is not, of substantive importance) would probably remain unchanged.
- The impact of bias may be substantial. If all relevant studies were included, the key finding (that the effect size is, or is not, of substantive importance) could change.

Critically, one should never assert that the adjusted value is the ‘Correct’ value (Borenstein, 2019; Carter, Schönbrodt, Gervais, & Hilgard, 2019; Duval & Tweedie, 2000b; McShane, Bockenholt, & Hansen, 2016; Peters, Sutton, Jones, Abrams, & Rushton, 2007; Vevea *et al.*, 2019). The reader should be skeptical of several increasingly popular methods for examining publication bias known as *p*-curve (Simonsohn, Nelson, & Simmons, 2014) and PET-PEESE (Stanley & Doucouliagos, 2014) which maintain that the publication-bias adjusted effect is the true effect.

### SOME IMPORTANT CAVEATS

The procedures outlined in this chapter allow us to assess the relationship between effect size and study size. If this relationship does exist, it could reflect publication bias but could also reflect the fact that the effect size is larger in smaller studies. We need to address one additional issue. We can only apply these procedures if certain

basic conditions are met. Unless these conditions are met, the idea of looking for a relationship is pointless (Ioannidis & Trikalinos, 2007; Sterne *et al.*, 2011).

- We need to have a reasonable number of studies. There is a consensus that we should use 10 studies as a minimum, but that many more studies would be needed in the presence of substantial heterogeneity (Higgins & Thomas, 2019; Ioannidis & Trikalinos, 2007; Sterne *et al.*, 2011).
- Even with a reasonable number of studies, statistical power to identify a relationship will often be low. Therefore, our failure to find evidence of asymmetry should not lead to a false sense of assurance.
- We need to have a reasonable amount of variation in the sample size. The procedures all look for a relationship between effect size and sample size. If all studies have approximately the same sample size, then (by definition) there can be no correlation between sample size and effect size. When all studies in the analysis have been performed by drug companies, there is a distinct possibility that all studies will have a similar sample size, since drug companies often use a standard sample size for a particular type of study.
- We need to have a reasonable amount of variation in the observed effect sizes. If all studies have approximately the same effect size, then (by definition) there can be no correlation between sample size and effect size.
- The methods may yield a very different picture depending on the index used in the analysis (e.g. risk difference versus risk ratio).
- There must be at least one study in the analysis that is statistically significant. If no studies are statistically significant, it makes no sense to suggest that our sample was biased by the preferential inclusion of statistically significant studies (Ioannidis & Trikalinos, 2007).

### **PROCEDURES DO NOT APPLY TO STUDIES OF PREVALENCE**

The idea that studies are more likely to be published when they are statistically significant only makes sense when the studies in the analysis test for statistical significance. It does not apply to studies that report the prevalence of a condition. When we report a prevalence, we simply report the prevalence. We do not test it to see whether it is significantly different from any specific value. The possibility that a study will not be published because it is not statistically significant simply does not apply when there is no test of significance.

### **THE MODEL FOR PUBLICATION BIAS IS SIMPLISTIC**

The procedures being discussed are all based on a model that assumes studies which are statistically significant are more likely to find their way into an analysis than studies which are not statistically significant. Research shows that these assumptions do mirror the true state of affairs, in general. However, it is important to recognize that this is a

simplistic view of the overall situation and may not apply in any given case (Bax & Moons, 2011).

For example, one could imagine a scenario where the *first* studies looking at the impact of a new intervention are more likely to be published if they are statistically significant. After that, studies that confirm the original findings might be *harder* to publish, since they (merely) confirm what we already know, while studies that are *not* statistically significant might be published more readily, since they challenge the current state of information.

Additionally, the basic idea that statistically significant studies are more likely to be published than nonsignificant studies varies by trends and by journals. Recent awareness of the problems with lack of replication may shift editorial priorities, and some journals will agree to publish any study provided that the protocol is of high value and is followed correctly.

In this chapter, we have focused exclusively on procedures that are in common use and are relatively simple to use. More advanced procedures are discussed by Bayarri (1988), DuMouchel (1988), Hedges (1988, 1992), Iyengar and Greenhouse (1988a, 1988b), Keith and Begg (1992), Laird, Patil, and Taillie (1988), Rao (1988), and Rosenthal and Rubin (1988).

## CONCLUDING REMARKS

Almost any meta-analysis where studies are pulled from the literature will be affected by publication bias. Fortunately, that does not necessarily invalidate the analysis. The key issue is not whether *any* bias exists, but rather *how much* of an impact this bias might have caused. In many cases, we will be able to say that while bias may have inflated the effect size, the basic conclusions of the analysis are robust.

While publication bias is often a problem, it would help to look at publication bias in the context of the entire meta-analysis procedure. Publication bias is only one factor that impacts the estimate of the mean effect size. For example, consider an analysis where our goal is to estimate the mean impact of a drug for all populations within a defined universe. For the analysis to yield an accurate estimate of this mean would require that:

- A. The studies that had actually been performed to address the impact of the drug are a random sample of all relevant populations, and of all variants on the intervention.
- B. The studies that were retrieved for our analysis are a random sample of all studies that had been performed.
- C. The outcome reported in each paper was either the outcome identified *a priori* or was randomly selected from all outcomes measured.

In fact, none of these is likely to be true. Our estimate of the mean will be affected by problems with (A), (B), and (C). When viewed in this context, it should be clear that the utility of a meta-analysis is not dependent primarily on (B). If publication bias is a problem, the analysis may still yield useful information. If publication bias is not a serious problem, we need to recognize that the estimates may still be affected by other factors.

## PUTTING IT ALL TOGETHER

Publication bias refers to the fact that studies which overestimate the impact of an intervention are more likely to be included in meta-analyses than studies which accurately estimate or underestimate the impact of that intervention. Researchers have developed methods that are intended to identify the presence of publication bias. While these methods are sometimes helpful, they have some serious limitations.

These methods look for a relationship between study size and effect size. If the effect size tends to go up as the sample size goes down, this could be evidence of publication bias. However, this could also reflect the fact that the effect size actually is larger in smaller studies for reasons that are unrelated to bias. For that reason, the relationship should be called a small-study effect rather than publication bias.

Since there is no way to distinguish between a small-study effect and publication bias, the methods which yield an adjusted effect size must be seen as a sensitivity analysis rather than seen as estimating the correct effect size. Using these methods, we can report that if we adjusted the effect to remove the bias, (a) the resulting effect would be essentially unchanged, or (b) the effect might change but the basic conclusion, that the treatment works (or not) would not be changed, or (c) the basic conclusion would be called into question. *In no case* should the results of a publication bias analysis be regarded as estimating the true effect.

### SUMMARY POINTS

- Publication bias exists when the studies included in the analysis differ systematically from all studies that should have been included. Typically, studies with larger than average effects are more likely to be published, and this can lead to an upward bias in the summary effect.
- Methods have been developed that look for a relationship between study size and effect size. This relationship should be called a small-study effect. This could be due to publication bias, but could also reflect the fact that the effect size actually is larger in smaller studies.
- Since we cannot determine whether the small-study effect reflects bias (on the one hand) or the fact that the effect size really is larger in small studies (on the other), the methods employed to assess publication bias cannot be used to determine the 'correct' effect size.
- The methods should be used to serve as a sensitivity analysis. We could report that if the relationship between study size and effect size really is due to bias, (a) the results would be essentially unchanged, or (b) the mean effect size might change but the basic conclusion would not be changed, or (c) the basic conclusion would be called into question.

## Further Reading

- Boutron, I., Page, M.J., Higgins, J.P.T., Altman, D.G., Lundh, A., & Hróbjartsson, A. (2019). Considering bias and conflicts of interest among the included studies. In J.P.T. Higgins, J. Thomas, M. Cumpston, T. Li, M.J. Page, & V.A. Welch (eds), *Cochrane Handbook for Systematic Reviews of Interventions* (2nd edition, pp. 177–204). Chichester, UK: John Wiley and Sons, Ltd.
- Chalmers, T.C., Frank, C.S., & Reitman, D. (1990). Minimizing the three states of publication bias. *JAMA* 263: 1392–1395.
- Dickersin, K., Chan, S., Chalmers, T.C., Sacks, H.S., & Smith, H. (1987). Publication bias in clinical trials. *Controlled Clinical Trials* 8: 348–353.
- Dickersin, K., Min, Y.L., & Meinert, C.L. (1992). Factors influencing publication of research results: follow-up of applications submitted to two institutional review boards. *JAMA* 267: 374–378.
- Hedges, L.V. (1984). Estimation of effect size under nonrandom sampling: the effects of censoring studies yielding statistically insignificant mean differences. *Journal of Educational Statistics* 9: 61–85.
- Hedges, L.V. (1989). Estimating the normal mean and variance under a selection model. In Gleser, L., Perlman, M.D., Press, S.J., Sampson, A.R. (eds), *Contributions to Probability and Statistics: Essays in Honor of Ingram Olkin* (pp. 447–458). New York, NY: Springer Verlag.
- Hunt, M.M. (1999). *The New Know-Nothings: The Political Foes of the Scientific Study of Human Nature*. New Brunswick, NJ, Transacation.
- International Committee of Medical Journal Editors. *Uniform Requirements for Manuscripts Submitted to Biomedical Journals: Writing and Editing for Biomedical Publication*. Updated October 2007. Available at [http://www.icmje.org/#clin\\_trials](http://www.icmje.org/#clin_trials).
- Ioannidis, J.P. (2007). Why most published research findings are false: author's reply to Goodman and Greenland. *PLoS Medicine* 4: e215.
- Lau, J., Ioannidis, J.P., Terrin, N., Schmid, C.H., & Olkin, I. (2006). The case of the misleading funnel plot. *BMJ* 333: 597–600.
- Lefebvre, C., Glanville, J., Briscoe, S., et al. (2019). Searching for and selecting studies. In J.P.T. Higgins, J. Thomas, M. Cumpston, T. Li, M.J. Page, & V.A. Welch (eds), *Cochrane Handbook for Systematic Reviews of Interventions* (2nd edition, pp. 67–108). Chichester, UK: John Wiley and Sons, Ltd.
- Rosenthal, R. (1979). The 'File drawer problem' and tolerance for null results. *Psychological Bulletin* 86: 638–641.
- Rothstein, H.R., Sutton, A.J., & Borenstein, M. (2005). *Publication Bias in Meta-analysis: Prevention, Assessment and Adjustments*. Chichester, UK: John Wiley & Sons, Ltd.
- Sterne, J.A., Egger, M. & Smith, G. D. (2001). Systematic reviews in healthcare: investigating and dealing with publication and other biases in meta-analysis. *BMJ* 323: 101–105.
- Sterne, J.A., Egger, M. (2001). Funnel plots for detecting bias in meta-analysis: guidelines on choice of axis. *Journal of Clinical Epidemiology* 54: 1046–1055.
- Sutton A.J., Duval S.J., Tweedie, R.L., Abrams, K.R., & Jones, D.R. (2000). Empirical assessment of effect of publication bias on meta-analyses. *BMJ* 320: 1574–1577.



## Issues Related to Effect Size



# Overview

The focus of this volume has been almost exclusively on meta-analysis of effect sizes (though some other approaches will be addressed in Chapter 41). We compute an effect size for each study, and it is these values that form the core of the analysis. We use them to assess dispersion across studies, to compute a summary effect, to compare effects across subgroups, and so on. In this part we offer some context for this approach.

Readers who are accustomed to working with significance tests and  $p$ -values in primary studies may wonder why we don't synthesize these  $p$ -values. This is addressed in Chapter 37.

Some may wonder why we compute an effect size for each study, rather than simply aggregating the summary data (for example, summing all the cell counts for a series of  $2 \times 2$  tables) and then computing an effect size for the final table. This is addressed in Chapter 38.

In this volume we have worked primarily with measures that assess the impact of treatments or the relationship between variables, but the same idea can be extended to other kinds of measures as well. In Chapter 39 we provide an overview of some other applications, such as estimating prevalence. We also outline some other approaches to meta-analysis, such as the use of individual participant data, and of Bayesian meta-analysis.



# Effect Sizes Rather than $p$ -Values

---

### Introduction

Relationship between  $p$ -values and effect sizes

The distinction is important

The  $p$ -value is often misinterpreted

Narrative reviews vs. meta-analyses

---

## INTRODUCTION

A central theme in this volume is the fact that we usually prefer to work with effect sizes, rather than  $p$ -values. Readers who are accustomed to working with significance tests and  $p$ -values in primary studies may wonder why we don't synthesize these  $p$ -values. The reason reflects a fundamental issue that applies both to primary studies and to meta-analysis, and is the subject of this chapter. Since narrative reviews typically work with  $p$ -values while meta-analyses typically work with effect sizes, the distinction between effect sizes and  $p$ -values also reflects a difference between narrative reviews and meta-analysis, and we address this as well.

Note that meta-analysis methods are available for working with  $p$ -values, and we will describe these in Chapter 41.

## RELATIONSHIP BETWEEN $p$ -VALUES AND EFFECT SIZES

There are two general approaches to data analysis, both in primary studies and in meta-analysis. One is significance testing, where the researcher poses the null hypothesis (for example, that the treatment effect is zero) and then attempts to disprove that hypothesis. The other is effect-size estimation, where the researcher estimates the magnitude of the effect size. Both approaches start with the same values but express them in different ways.

Suppose a study of independent groups reports a standardized mean difference of 0.50 with a standard error of 0.20. To perform a significance test we compute

$$Z = \frac{d}{SE_d}, \quad (37.1)$$

and then compare the observed Z-value to the Z-value required for statistical significance. In this example

$$Z = \frac{0.50}{0.20} = 2.5000.$$

We then compare the observed Z-value of 2.50 with the Z-value of 1.96 (corresponding to the two-tailed alpha of 0.05) and conclude that the effect is statistically significant.

The complementary approach is to report the effect size (which reflects the magnitude of the effect) and its confidence interval (which reflects the precision of the estimate). In this example the effect size is  $d = 0.50$ . The confidence interval is given by

$$LL_d = d - 1.96 \times SE_d \quad (37.2)$$

and

$$UL_d = d + 1.96 \times SE_d \quad (37.3)$$

where 1.96 corresponds to the two-tailed 0.05 significance value. Concretely, the lower limit and upper limit for  $d$  are

$$LL_d = 0.50 - 1.96 \times 0.20 = 0.1080$$

and

$$UL_d = 0.50 + 1.96 \times 0.20 = 0.8920.$$

The two approaches are consistent with each other in that the  $p$ -value for a test of the null hypothesis will fall under 0.05 if and only if the 95% confidence interval for the effect size excludes the null value. We can see this in Figure 37.1. The top line reflects the effect size approach, with a line representing the confidence interval extending from the effect size toward the null value. The bottom line reflects the significance testing approach, with a line representing the nonsignificance region extending from the null value toward the effect size. The line is the same length in either case ( $Z \times SE_d$ ), which means that the top line will include the null value if and only if the bottom line includes the effect size.

Note. There may be small differences in the length of the line but the difference is generally trivial.

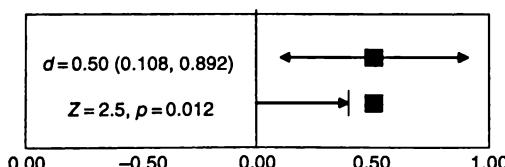


Figure 37.1 Estimating the effect size versus testing the null hypothesis.

## THE DISTINCTION IS IMPORTANT

Since the two approaches are consistent with each other, the decision to use one or the other is largely a matter of choice. However, it is an important choice, because by focusing on different questions (the viability of the null hypothesis versus the magnitude of the effect) the two approaches shift the attention of the researcher in important ways. While researchers in many fields tend to favor significance tests, there are good reasons to focus on effect sizes instead.

First, the effect size is what we need to decide on a course of action. If a clinician or patient needs to make a decision about whether or not to employ a treatment, they want to know if the treatment reduces the risk of death by 5% or 10% or 20%, and this is the information carried by the effect size. Similarly, if we are thinking of implementing an intervention to increase the test scores of students, or to reduce the number of incarcerations among at-risk juveniles, or to increase the survival time for patients with pancreatic cancer, the question we ask is about the magnitude of the effect. The *p*-value tells us only that the effect may be (or is probably not) zero.

Second, the *p*-value is often misinterpreted. Because researchers *care about* the effect size, they tend to take whatever information they have and press it into service as an indicator of effect size. A statistically significant *p*-value is assumed to reflect a clinically important effect, and a nonsignificant *p*-value is assumed to reflect a trivial (or zero) effect. However, these interpretations are not necessarily correct. The problem with using the *p*-value as a surrogate for effect size is that the *p*-value incorporates information about both the size of the effect and also the size of the sample (or the precision with which the effect is estimated). While a significant *p*-value *may* reflect a large effect size, it could also reflect a small effect size that had been measured in a large study. Similarly, while a nonsignificant *p*-value *may* reflect a small effect size, it could also reflect a large effect size that had been measured in a small study.

For example, Figure 37.2 is a plot of two studies. Study A has a *p*-value of 0.119 while Study B has a *p*-value of < 0.001, which might suggest that the effect is stronger in Study B. In fact, the effect size is the same (0.50) in both studies. The difference in *p*-values reflects a difference in sample size (40 in Study A versus 200 in Study B), not a difference in effect size.

By contrast, when we work with the effect size we focus on the question of interest, which is to estimate the magnitude of the effect. We report the effect size as 0.50,

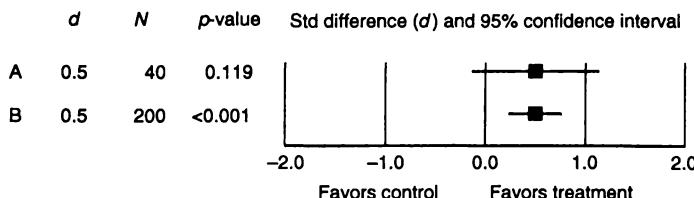


Figure 37.2 The *p*-value is a poor surrogate for effect size.

and (as a separate matter) the precision ( $-0.129$  to  $1.129$  for Study *A*, and  $0.219$  to  $0.781$  for Study *B*). Additionally, this approach avoids the mistakes outlined in the previous paragraph. Because we work with the effect sizes directly we avoid the problem of interpreting nonsignificant  $p$ -values to indicate the absence of an effect (or of interpreting significant  $p$ -values to indicate a large effect).

### THE $p$ -VALUE IS OFTEN MISINTERPRETED

While we would argue that researchers should shift their focus to effect sizes even when working entirely with primary studies, the shift is *absolutely critical* when our goal is to synthesize data from multiple studies. A narrative reviewer who works with  $p$ -values (or with reports that were based on  $p$ -values) and uses these as the basis for a synthesis, is facing an impossible task. Where people tend to misinterpret a single  $p$ -value, the problem is much worse when they need to compare a series of  $p$ -values. Consider the following three examples.

Suppose we are told that four studies reported  $p$ -values of  $0.28$ ,  $0.28$ ,  $0.28$ , and  $0.003$ . A reviewer working with these  $p$ -values might assume that the effect was larger in the last study. The studies are shown in Figure 37.3. The effect size is the same in all the studies, and the  $p$ -values differ only because the sample size was larger in the last study.

Suppose we are told that three studies each reported a  $p$ -value of  $0.012$ . Many would assume that the treatment effect is consistent across studies. The studies are shown in Figure 37.4. Study *A* has a large effect (and poor precision), Study *B* has a moderate effect (and modest precision), while Study *C* has a small effect (and excellent precision).

Suppose we are told that studies *A* and *B* reported  $p$ -values of  $0.057$ , and  $0.035$ . Many would assume that the effect size was higher in Study *B*. The studies are shown in Figure 37.5, and it turns out that the effect was *weaker* in study *B*.

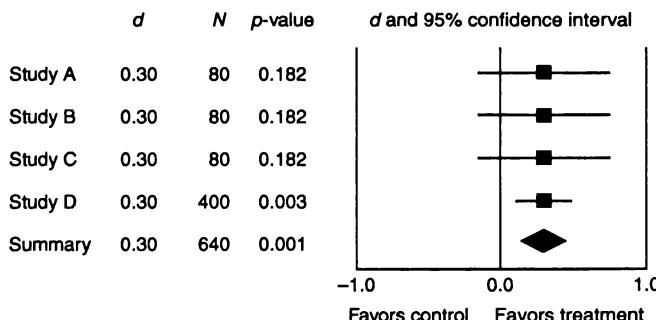


Figure 37.3 Studies where  $p$ -values differ but effect sizes is the same.

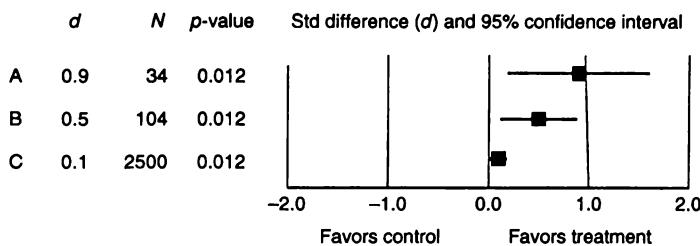


Figure 37.4 Studies where *p*-values are the same but effect sizes differ.

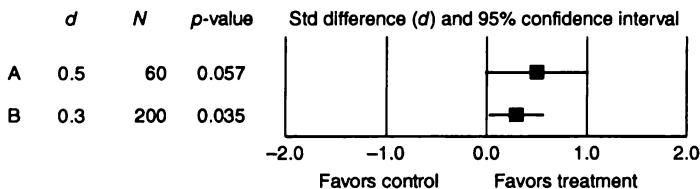


Figure 37.5 Studies where the more significant *p*-value corresponds to weaker effect size.

## NARRATIVE REVIEWS VS. META-ANALYSES

The narrative review typically works with *p*-values (or with conclusions that are based on *p*-values), and therefore lends itself to these kinds of mistakes. *p*-values that differ are assumed to reflect different effect sizes but may not (Figure 37.3), *p*-values that are the same are assumed to reflect similar effect sizes but may not (Figure 37.4), and a more significant *p*-value is assumed to reflect a larger effect size when it may actually be based on a smaller effect size (Figure 37.5). By contrast, the meta-analysis works with effect sizes. As such it not only focuses on the question of interest (what is the size of the effect) but allows us to compare the effect size from study to study.

There is an additional difference between meta-analysis and narrative reviews. Where the narrative review treats each item of information as discrete (and the synthesis takes place in the reviewer's head), meta-analysis incorporates all of the effect sizes in a single statistical analysis. Even if the *p*-value did move in tandem with the effect size (for example, if all studies in the analysis had the same sample size), the narrative review provides no mechanism for assessing the dispersion in effect size from one study to the next. By contrast, meta-analysis usually works directly with the effect size (which is separated from the sample size), and uses established statistical techniques to isolate and quantify the true dispersion.

**SUMMARY POINTS**

- To synthesize data from a series of studies we need to work with the effect size rather than the *p*-value from each study and we need to incorporate all of the effect sizes in a single analysis, rather than working with discrete results from the separate studies.
- The narrative review fails on both counts. It works with the *p*-values from the primary studies (either directly, or because most studies base their results section and discussion on the *p*-values). And, it tries to perform a synthesis working with a series of discrete results.
- By contrast, meta-analysis has developed methods to meet these two goals. It works with the effect sizes, and incorporates all of these in a single analysis.

---

## CHAPTER 38

---

# Simpson's Paradox

---

### Introduction

Circumcision and risk of HIV infection

An example of the paradox

---

## INTRODUCTION

To compute the summary effect in a meta-analysis we compute an effect size for each study and then combine these effect sizes, rather than pooling the data directly. For example, if we start with  $2 \times 2$  tables we compute an odds ratio for each table and then combine these odds ratios. We *do not* pool the cell counts across tables to create a pooled  $2 \times 2$  table and then compute the odds ratio for this table.

This approach allows us to study the dispersion of effects before proceeding to the summary effect. For a random-effects model this approach also allows us to incorporate the between-studies dispersion into the weights.

There is one additional reason for using this approach, and that reason is the subject of this chapter. The reason is to ensure that each effect size is based on the comparison of a group with *its own* control group, and thus avoid a problem known as Simpson's paradox. In some cases, particularly when we are working with observational studies, this is a critically important feature.

To illustrate this point we present a review published by Van Howe (1999) which concluded that circumcision is associated with increased risk of HIV. Matthias Egger presented a critique of this review at the Cochrane Colloquium in Cape Town in 1999, and (with O'Farrell) also published a commentary in the *International Journal of STD and AIDS* (O'Farrell and Egger, 2000). What follows draws heavily on this work.

## CIRCUMCISION AND RISK OF HIV INFECTION

Van Howe (1999) published a review article in the *International Journal of STD and AIDS* that looked at the relationship between circumcision and HIV infection in Africa. The article was based on data from 33 studies, which Van Howe classified into one of

three groups based on the populations studied. The *High-risk* populations included long-distance truck drivers and patients recruited at STD clinics. The *Partner* studies looked at HIV infection in men whose partner was HIV positive. The *Random population* surveys did not target specific groups. The prevalence of HIV in the men in these three groups was 25%, 11% and 9%, respectively (see Table 38.1).

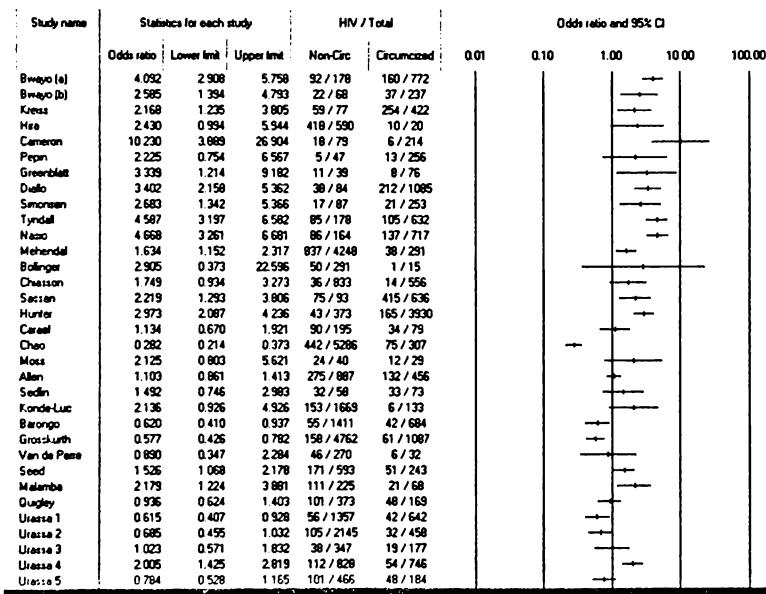
**Table 38.1** HIV as function of circumcision (by subgroup).

Sample	Studies	N	Prevalence HIV
High risk	15	13,238	25%
Partner	7	13,515	11%
Random	11	17,267	9%
<b>Total</b>	<b>33</b>	<b>44,020</b>	

The data for all 33 studies are shown in Table 38.2. An odds ratio less than 1.0 means that circumcision is associated with *higher* risk of HIV, and an odds ratio greater than 1.0 means that circumcision is associated with *lower* risk of HIV.

In the first study, Bwayo (a), the risk of HIV is 52% (92/178) for noncircumcised versus 21% (160/772) for circumcised, yielding an odds ratio of 4.09 (lower risk for circumcised), and so on for the remaining studies. Of the 33 studies, 8 associated circumcision with higher risk (an odds ratio of less than 1.0), and 25 with lower risk (an odds ratio greater than 1.0).

The same data are shown graphically in Figure 38.1.



**Figure 38.1** Circumcision and HIV. Odds Ratio >1 indicates circumcision is associated with lower risk of HIV.

## AN EXAMPLE OF THE PARADOX

Van Howe summed the counts in each cell across the 33 studies (see the line labeled *Total* in Table 38.2) to yield Table 38.3. Using this table he computed the risk of HIV as 14% (3962/28341) for noncircumcised versus 15% (2312/15679) for circumcised, and the odds ratio as 0.94 (95% confidence interval 0.89 to 0.99) with *p*-value under 0.05. Van Howe concludes that 'When the raw data are combined, a man with a circumcised penis is at *greater risk* of acquiring and transmitting HIV than a man with a non-circumcised penis.'

**Table 38.2** HIV as function of circumcision – by study.

Study name	Noncircumcised		Circumcised		Prevalence		Odds ratio
	HIV+	Total	HIV+	Total	Noncircumcised	Circumcised	
<b>High risk</b>							
Bwayo (a)	92	178	160	772	0.52	0.21	4.09
Bwayo (b)	22	68	37	237	0.32	0.16	2.59
Kreiss	59	77	254	422	0.77	0.60	2.17
Hira	418	590	10	20	0.71	0.50	2.43
Cameron	18	79	6	214	0.23	0.03	10.23
Pepin	5	47	13	256	0.11	0.05	2.23
Greenblatt	11	39	8	76	0.28	0.11	3.34
Diallo	38	84	212	1085	0.45	0.20	3.40
Simonsen	17	87	21	253	0.20	0.08	2.68
Tyndall	85	178	105	632	0.48	0.17	4.59
Nasio	86	164	137	717	0.52	0.19	4.67
Mehendal	837	4248	38	291	0.20	0.13	1.63
Bollinger	50	291	1	15	0.17	0.07	2.90
Chiasson	36	833	14	556	0.04	0.03	1.75
Sassan	75	93	415	636	0.81	0.65	2.22
<b>Partner</b>							
Hunter	43	373	165	3930	0.12	0.04	2.97
Carael	90	195	34	79	0.46	0.43	1.13
Chao	442	5286	75	307	0.08	0.24	0.28
Moss	24	40	12	29	0.60	0.41	2.13
Allen	275	887	132	456	0.31	0.29	1.10
Sedlin	32	58	33	73	0.55	0.45	1.49
Konde-Luc	153	1669	6	133	0.09	0.05	2.14
<b>Random population</b>							
Barongo	55	1411	42	684	0.04	0.06	0.62
Grosskurth	158	4762	61	1087	0.03	0.06	0.58
Van de Perre	46	270	6	32	0.17	0.19	0.89
Seed	171	593	51	243	0.29	0.21	1.53
Malamba	111	225	21	68	0.49	0.31	2.18
Quigley	101	373	48	169	0.27	0.28	0.94
Urassa 1	56	1357	42	642	0.04	0.07	0.61
Urassa 2	105	2145	32	458	0.05	0.07	0.69
Urassa 3	38	347	19	177	0.11	0.11	1.02
Urassa 4	112	828	54	746	0.14	0.07	2.00
Urassa 5	101	466	48	184	0.22	0.26	0.78
<b>Total</b>	3962	28341	2312	15679	0.14	0.15	0.94

**Table 38.3** HIV as a function of circumcision – full population.

	HIV Positive	HIV Negative	Total	%HIV
Noncircumcised	3962	24379	28341	14%
Circumcised	2312	13367	15679	15%

This conclusion seems to be at odds with the full table of data, where the preponderance of studies (25 of the 33) had odds ratios greater than 1.0, meaning that circumcision was associated with *lower risk* of HIV. This anomaly is also evident on the plot, where 25 studies line up on the right-side of 1.0, with only 8 on the left. Also, the 8 studies on the left are not large enough to have pulled the effect so far to the left. Therefore, it seems counter-intuitive that the summary effect should fall to the left.

In fact, this intuition is correct. When the data are analyzed using standard meta-analysis techniques, circumcision is associated with a *lower risk* of HIV. Under the random-effects model the odds ratio is 1.67 with a 95% confidence interval of 1.25 to 2.24, Z of 3.42 and *p*-value of 0.001. The fixed-effect model is not appropriate here since the effects are clearly heterogeneous ( $Q = 419$ ,  $df = 32$ ,  $p < 0.0001$ ) but the fixed-effect model would lead to the same conclusion with an odds ratio of 1.40, 95% confidence interval of 1.29 to 1.51, Z = 8.43, and a *p*-value of < 0.0001. These findings are consistent with the visual impression of the data, and contradict Van Howe's conclusions.

The reason that Van Howe reported a *higher risk* for circumcision is that rather than compute an effect size for each study and then pool these effects, he pooled the  $2 \times 2$  tables for each study and then used the combined table to compute the odds ratio. The reason that this matters is as follows.

Recall that the 33 studies included three kinds of populations, *High risk*, *Partners*, and *Random*, with prevalence rates of 25%, 11% and 9%, respectively. It turns out that a disproportionate number of the circumcised patients were from the *High risk* studies. Among circumcised persons 39% came from the *high risk* studies while only 29% came from the *Random* studies. Among the noncircumcised, by contrast, 25% came from the high risk studies while 45% came from the *Random* studies. The proportion of noncircumcised and circumcised subjects drawn from each kind of study is shown in Table 38.4.

In this light, the labels in Table 38.3 might better reflect the facts if we modified them as shown in Table 38.5. Rather than *Noncircumcised* vs. *Circumcised*, the rows are labeled *Noncircumcised, low HIV prevalence population* vs. *Circumcised, high HIV prevalence population*.

**Table 38.4** HIV as a function of circumcision – by risk group.

	High risk	Partner	Random
Noncircumcised	25%	30%	45%
Circumcised	39%	32%	29%

**Table 38.5** HIV as a function of circumcision/risk group – full population.

	HIV Positive	HIV Negative	Total	%HIV
Noncircumcised, low HIV prevalence population	3962	24379	28341	14%
Circumcised, high HIV prevalence population	2312	13367	15679	15%

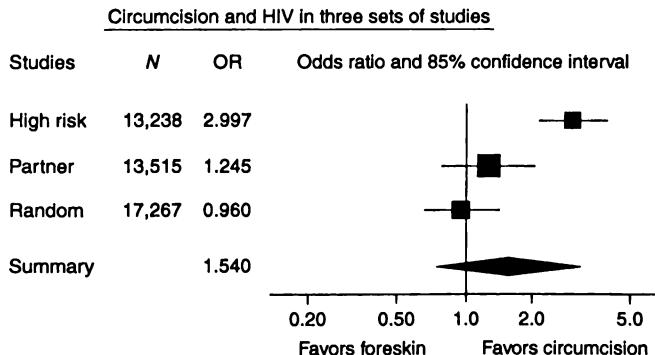
In sum, the problem with Van Howe's analysis is that by summing the raw data he introduced confounding. That is, the noncircumcised men cannot legitimately be compared with the circumcised population because they differ in other ways. The increased risk of HIV that appeared to be associated with circumcision was also associated with an array of high-risk behaviors (and more likely due to those than to the circumcision). By contrast, when we compute an effect size for each study we control for these confounds (at least to the extent possible when working with observational studies).

Egger also showed what happens if we compute the odds ratio *within* each of the three groups that Van Howe had delineated. Within the high risk group, circumcision is associated with a decreased risk (by about two-thirds) of HIV. Within each of the other two groups of studies the relationship between circumcision and HIV (if any) is small (see Figure 38.2).

Egger used this pattern of effects to formulate several hypotheses about the mechanisms by which circumcision might have an impact on the risk of HIV infection, and suggested that these be tested in future studies. This is an elegant and appropriate use of meta-analysis, to direct the formulation and conduct of future studies.

Egger also noted that given the heterogeneity of effects it would be a mistake to draw any conclusions from the analysis at this point. As Egger said in concluding his presentation at Capetown, 'I am not suggesting that you all run right out and get circumcised.'

The problem addressed in this chapter has been known to statisticians for a long time, and is generally called Simpson's paradox. The term paradox refers to the fact

**Figure 38.2** HIV as function of circumcision – in three sets of studies.

that one group can do better in *every one* of the included studies, but still do worse when the raw data are pooled. The problem is not limited to studies that use proportions, but can exist also in studies that use means or other indices. The problem exists only when the base rate (or mean) varies from study to study and the proportion of participants from each group varies as well. For this reason, the problem is generally limited to observational studies, although it can exist in randomized trials when allocation ratios vary from study to study.

For readers who choose to savor these papers in the original, note that the Van Howe paper seems to include a number of typographical errors which we corrected for this chapter (in particular, the Van Howe abstract reverses the direction of the odds ratio from the one used in the body of paper). Note also that Egger raises many questions about Van Howe's methodology in addition to the ones outlined here.

### SUMMARY POINTS

- In a meta-analysis we compute an effect size for each study and combine these effect sizes, rather than combining the summary data and then computing an effect size for the combined data.
- This allows us to determine whether or not the effects are consistent across studies.
- This also ensures that each study serves as its own control, which minimizes the potential impact of confounders.
- If we were to pool data across studies and then compute the effect size from the pooled data, we may get the wrong answer, due to Simpson's paradox.

### Further Reading

- Chudnovsky, A & Niederberger, C.S. (2007). The foreskin dilemma: To cut or not to cut. *J Androl* 28: 5–7.
- O'Farrell, N., & Egger, M. (2000). Circumcision in men and the prevention of HIV infection: a 'meta-analysis' revisited. *Int J STD AIDS* 11: 137–142.
- Siegfried, N., Muller, M., Volmink, J., et al. (2003). Male circumcision for prevention of heterosexual acquisition of HIV in men. Cochrane Database of Systematic Reviews Issue 3. Art. No.: CD003362. DOI: 10.1002/14651858.CD003362.
- Siegfried, N., Muller, M., Deeks, J., et al. (2005). HIV and male circumcision-a systematic review with assessment of the quality of studies. *Lancet Infect Dis* 5: 165–173.
- Van Howe, R.S. (1999). Circumcision and HIV infection: review of the literature and meta analysis. *Int J STD AIDS* 10(1): 8–16.
- Van Howe R.S., Cold, C.J., Storms, M.R., et al. (2000). Male circumcision and HIV prevention. *BMJ* 321:1467–1468; 1469.
- Van Howe R.S., Svobodo, J.S., & Hodges, M. (2005). HIV infection and circumcision: cutting through the hyperbole. *Journal of the Royal Society for the Promotion of Health* 125: 259–265.

# Generality of the Basic Inverse-Variance Method

---

Introduction

Other effect sizes

Other methods for estimating effect sizes

Individual participant data meta-analyses

Bayesian approaches

---

## INTRODUCTION

The basic idea of meta-analysis is to compute an effect size from each of several studies, and to calculate a weighted average of these effect size estimates. We have shown how the weights can reflect either a fixed-effect assumption or a random-effects assumption. We have also seen how we can examine differences in effect sizes across studies, using subgroup analyses and meta-regression.

In all of our examples so far, we have applied these methods to *comparative studies* or studies of *association*. For example, we might compare the outcomes of people randomized to different treatments in a clinical trial by looking at the difference in means or the ratio of risks. Or, we could use the same techniques to compare the scores in two existing groups, such as males and females.

Precisely the same approach to meta-analysis can be used to work with other kinds of studies and other kinds of data. We will use the term *point estimate* as a more generic term than effect size estimate, to reflect the fact that we are not necessarily looking at the effect of one thing on another. The only requirements for meta-analysis to be possible are that:

- The point estimate can be expressed as a single number.
- We can compute a variance for the point estimate.

## OTHER EFFECT SIZES

We now provide some examples of situations in which these requirements are met and where meta-analysis can therefore be used to combine findings across studies. We do not attempt to provide detailed methodology for any of these. Rather, our aim is to provide a flavor of the range of applications for the basic methods. Many of the relevant formulas are provided in the companion volume, *Computing Effect Sizes for Meta-Analysis* (Borenstein, Hedges, Higgins and Rothstein, in preparation, John Wiley & Sons, Ltd).

### Single descriptive statistics

A meta-analysis can be used to combine single descriptive statistics. For example, a sample of continuous measures might be summarized simply by its mean, and a meta-analysis might be performed to synthesize means across samples (or studies).

Similarly, a sample of binary outcomes might be summarized simply by the proportion of successes, the risk of an event, or the prevalence of a condition. For example, a meta-analysis was used to assess the prevalence of food allergies in the general population (Rona *et al.*, 2007). Prevalence was estimated from each study, along with its variance, and combined in a meta-analysis using standard techniques. However, as is often the case for single-group studies, very substantial heterogeneity was observed. For instance, while the average prevalence of self-reported milk allergy was a little over 3% in a meta-analysis, with a tight 95% confidence interval from around 2.5% to 4%, the results from the individual studies ranged from 1.2% to 17%.

### Physical constants

Meta-analysis methods have a long history in the field of physics. For example, Raymond Birge published a paper on methods for combining estimates of physical constants from different experiments in 1932. In 1941, he summarized numerous attempts that had been made to measure the speed of light in a vacuum, and calculated weighted averages of their findings (Birge, 1941). A similar approach has been used for other important constants. These are essentially meta-analyses of the multiple experiments, although the term meta-analysis was not introduced until 1976.

### Two-group studies with other types of data

Although we have described in some detail different indices for comparing two groups when the outcomes are continuous or binary, we should make it clear that other types of data may be encountered. Three particular different types of data are as follows.

Ordinal data arise when each individual is assigned to one of *three or more* categories, and these categories have a logical ordering. For example, the symptoms of a condition after a period of treatment might be assessed as 'mild', 'moderate' or 'severe'. A single effect size can be computed from studies with ordinal data for use

in meta-analysis. For example, under assumptions of a proportional odds model, an odds ratio is available (this is not computed in the same way as in Chapter 5).

Time-to-event (survival) data are used when we know how long each person was followed, and the outcome (either the event occurred or the follow-up period ended). Typically we compute a hazard ratio to reflect the treatment effect in each study and use meta-analysis to synthesize these hazard ratios across studies.

Rate data arise when the data includes both a period of observation and the number of times a specific event occurs during this period. This type of data is used to determine whether the incidence of a (repeatable) event is lower in one group than the other. In particular, it allows for individuals to have more than one event.

Methods are available for comparing two groups when any of these types of data are encountered. As noted above, the sorts of indices that are most commonly encountered are odds ratios for ordinal data, hazard ratios for time-to-event data and rate ratios for rate data. The meta-analysis would proceed in the same way, in each of these particular cases working in the logarithmic scale because the indices are ratios rather than differences.

### Three-group studies

We provide one example of how basic meta-analysis methods have been applied to studies with three groups, in the field of genetic epidemiology. There is considerable interest in understanding how our genetic make-up affects our risk of developing diseases such as cancer and cardiovascular disease. Many studies are therefore investigating these relationships by measuring genetic variants and looking at the association between these variants and disease. Most genetic variants have two versions, called alleles. Crudely put, one of these alleles may usually be considered the *normal* variant, with the other having arisen as a mutation at some point during human evolution. For each genetic variant, we receive one copy from our father and one from our mother. We therefore end up with one of three possible combinations: two normal alleles, two mutant alleles or one normal allele and one mutant allele.

A simple genetic association study divides individuals into three groups according to their genetic variants, and cross-tabulates these groups against disease status. For example, one study by Lacasafía-Navarro and colleagues (2006) measured a specific variant in the methylenetetrahydrofolate reductase (MTHFR) gene, and compared the distributions of people having different combinations of alleles with their gastric cancer status. The two alleles for this specific variant (at position 677 in the gene) are termed 677C and 677T. The results are shown in Table 39.1.

There are several ways in which this study can be included in a basic meta-analysis. However, comparisons among three groups cannot be reduced to a single number without making assumptions or ignoring some of the information. A published meta-analysis by Boccia (2008) and colleagues created a  $2 \times 2$  table by excluding the middle (normal, mutation) group. Thus, they compared people with two mutation alleles with people with two normal alleles. This simple approach imposes no assumptions about the way in which the genetic variants determine risk, but does not make use of

**Table 39.1** Simple example of a genetic association study.

Genotype	Interpretation	Cases	Controls	Total
677C, 677C	normal, normal	56	144	200
677C, 677T	normal, mutation	85	179	264
677T, 677T	mutation, mutation	60	104	164
Total		201	427	628

all of the data. To use all of the data in a standard meta-analysis, we have to make assumptions. Three common assumptions are:

- A *dominant* genetic model, in which we assume that at least one mutation is sufficient to alter risk of disease (and it does not matter whether one or two mutation alleles is present).
- A *recessive* model, in which we assume both mutation alleles are required to alter risk of disease (and having one is no different from having none).
- An *additive* model, in which we assume there is a similar change in risk for each additional mutation allele.

To perform a meta-analysis assuming a dominant or a recessive model, the  $3 \times 2$  table from each study is collapsed into a  $2 \times 2$  table from each study, by combining either the second two rows (for a dominant model) or the first two rows (recessive model) in the table above. To perform an additive model, alternative methods need to be used for each study such as logistic regression. In every case, however, a log odds ratio and its variance is obtained for each study, and the meta-analysis proceeds in exactly the same way as we have described earlier in the book.

### Regression coefficients

A final example of the broad range of types of study amenable to meta-analysis is the combination of regression coefficients, or beta weights. This is the synthesis of results from several regression analyses, and should not to be confused with meta-regression. For example, Sirmans and colleagues (2006) were interested in the relationship between house price and various characteristics of the house such as its square footage, age, number of bedrooms and presence of a swimming pool. Tackling each of these separately, they compiled studies that have looked at the relationship (e.g. between square footage and house price) and obtained regression coefficients to characterize the relationship. These could then be combined in a meta-analysis, using the variances of the coefficients in the weights.

However, the main aim of these authors was not simply to obtain an overall regression coefficient, but to examine whether the regression coefficients depend on other study-level characteristics, such as the location and timing of the study, and whether or not the regression analysis controlled for other house characteristics (e.g. number of bedrooms). The primary analysis is therefore a meta-regression of regression coefficients.

## OTHER METHODS FOR ESTIMATING EFFECT SIZES

In the previous section we discussed how the basic meta-analysis method of computing inverse-variance weighted averages can be applied to diverse types of studies and to diverse types of data. In a fixed-effect meta-analysis, the weighted average provides a summary estimate of a common effect size. In a random-effects meta-analysis, the weighted average (using revised weights) provides an estimate of the mean effect size across studies. In a random-effects meta-analysis we also compute an estimate of the standard deviation of effect sizes across studies.

The inverse-variance weighted average approach is not the only way to perform a meta-analysis. Both fixed-effect meta-analyses and random-effects meta-analyses can be undertaken using a variety of other statistical methods. The need to consider other methods typically arises either to refine the analysis of data within a study, or to refine the analysis of variation across studies (or, quite frequently, for both).

### Refined methods tailored to the type of data within a study

First, we may wish to analyze the data *within a study* using methods specific to the type of data we have. For instance, if we have  $2 \times 2$  tables from binary data, the methods we have described involve the computation of an effect size for each study (such as a log odds ratio or a risk difference) and its variance. For most purposes this is appropriate, but when studies are small, or when events are rare, other methods can have better statistical properties. To understand why this might be, note that the basic inverse-variance method assumes that the variance from each study is truly the variance of that study's effect size estimate. However, in reality the variance is only an estimate of the true variance. For large studies, the estimate is close to the true variance and the assumption is not problematic. For small studies, the variance may not be well estimated.

Methods based directly on  $2 \times 2$  tables for binary data do not require us to estimate the variance. Some examples are the Mantel-Haenszel method and the one-step method commonly referred to as the Peto method. We discuss these methods in Chapter 42.

Another situation in which we might prefer to analyse a study using methods specific to the type of data is when we have the original data from each study, often referred to as *individual participant data*. We could compute a standardized mean difference, or an odds ratio, or some other effect size, from these data and perform an inverse-variance meta-analysis. However, we might alternatively wish to perform the meta-analysis directly on the complete original dataset. We discuss the potential advantages of this approach in the next section.

### Refined methods for analysing between-study variation

The second reason for considering other methods for meta-analysis is that we may wish to combine results *across studies* using more sophisticated statistical techniques.

In particular, in random-effects meta-analysis, we estimate the between-study standard deviation,  $\tau$ . When the number of studies is small, this can be estimated with considerable error, as we describe in Chapters 17 and 45. Methods are available that allow uncertainty in  $\tau$  to be taken into account. One such possibility is a Bayesian approach to the meta-analysis, which we discuss at the end of this chapter. A further consideration in the random-effects inverse-variance method is the assumption that the true effects in different studies follow a normal distribution. It is usually very difficult to assess whether this assumption is reasonable. Methods are available that assume other distributions, and even methods that allow the data to determine the shape of the distribution. These are advanced methods that so far have mostly been considered only by statistical methods researchers.

### INDIVIDUAL PARTICIPANT DATA META-ANALYSES

When the meta-analyst has access to all of the original data from each study, the meta-analysis may be referred to as an *individual participant data* (or *individual patient data*) meta-analysis. This usually involves collaboration with the authors of the original studies included in the meta-analysis. There are many advantages to individual participant data (IPD) meta-analysis over literature-based meta-analysis (or summary data meta-analysis), which are summarized by Stewart and Tierney (2002). These include:

- Being able to perform consistent data checking and (if necessary) data cleaning.
- Having available a complete and up-to-date dataset on which to base analyses.
- Being able to perform a wide variety of statistical analyses in the same way in every study.
- Being able to examine the effects of participant-level covariates.
- Further benefits of having direct contact with study authors, for example in collating descriptive information about the studies, interpreting results and identifying further studies.

With access to individual participant data, the range of possible analysis methods is substantial. Methods can be broadly categorized as methods that analyze each study separately and then combine effect sizes using standard meta-analysis techniques, and (on the other hand) methods that analyze all of the data in one go.

#### Applying standard meta-analysis methods to individual participant data

A common approach is to analyze each study in a consistent way and to perform an inverse-variance meta-analysis on the resulting effect size estimates and their variances. For example, a mean difference could be computed from each study, and these combined using the standard methods.

Having access to IPD, however, allows for consistent, and even complex, analyses of the data from each study. For instance, a problem sometimes encountered in summary data meta-analysis is that studies provide effect size estimates that are adjusted for different sets of covariates. For example, one study might adjust for age and sex; and

another might adjust for age, sex and smoking behavior. With access to the raw data, the meta-analyst can adjust for the same covariates in every study.

Another common limitation of summary data is that studies present basic results in such different ways that a common effect size cannot readily be computed for each study. This is particularly the case for time-to-event data. Because a pair of observations is collected on each individual (the length of observation and whether the event occurs at the end of this period), time-to-event data cannot conveniently be reduced to simple summaries such as a  $2 \times 2$  table. Thus, only results of *analyses* tend to be presented, which may differ across studies. Many of the existing IPD meta-analyses in the medical area address time-to-event outcomes, since a common method of analysis can be applied to each study prior to the meta-analysis.

### Analyzing individual patient data in a single analysis

Given IPD for a series of primary studies, any method that could be used to analyze the individual studies can be used to analyze the complete data set. The key principle underlying the analysis of an IPD dataset is that the individual identities of the studies are respected. In this way we avoid the problem of Simpson's paradox discussed in Chapter 38. In statistical terminology, we say that the analysis is stratified by study. This is often achieved in practice by including a dummy covariate for each study.

An advantage of analyzing all IPD together (rather than analyzing each study separately and then synthesizing the effect sizes) is that information may be borrowed from one study to another. For example, suppose we are interested in whether the effect of a weight-loss intervention depends on age, and we have several small trials of the intervention, each performed in a similar population with mixed ages. We cannot use meta-regression for such a question because age is a participant-level rather than a study-level covariate (the mean age will be roughly the same for every study). If we have IPD from the trials, including each participant's age, then we can obtain a powerful analysis of the effect of age on intervention by analyzing all of the data at once, providing we stratify by study. We can decide to borrow information from one study to another to further increase the power. For example, we could assume that the standard deviation for weight losses is the same in every study.

In fact, every one of the meta-analysis models we have discussed up to this point in the book can be performed in a single analysis of IPD, by making different assumptions in the statistical model. However, some of the random-effects models are difficult to implement. Therefore, the simpler approach of analyzing each study individually, and then synthesizing the effect sizes in a meta-analysis, remains the most popular in practice.

### BAYESIAN APPROACHES

The methods we describe in this book are classical, or frequentist, methods for statistics. They revolve around estimating unknown parameters along with confidence intervals, and performing statistical tests in order to determine the extent to which the results are compatible with a null hypothesis (the *p*-value).

An alternative approach to statistics is the Bayesian approach. This stems from a different philosophy of probability, and in particular from an interpretation of probability as an uncertainty (or a belief) rather than a frequency. Bayesian statistics attaches probability distributions to the parameters of interest. The main parameters of interest in a meta-analysis are the overall mean effect size and, in a random-effects meta-analysis, the standard deviation ( $\tau$ ) of true effect sizes across studies. We might, for example, represent our uncertainty about an overall mean log odds ratio by attaching a normal distribution centered on our *best guess* and with tails describing how confident we are about it.

A Bayesian analysis starts by attaching a prior probability distribution to each unknown quantity. This describes *a priori* uncertainty (or belief) about the quantity before seeing the data. In many cases, the prior distribution is used to express ignorance (e.g. as a flat distribution). The Bayesian analysis itself combines the prior distribution with the data, turning it into a posterior probability distribution for the unknown quantity. When the prior distribution represents prior ignorance, the posterior distribution is simply a summary of what the data tell us about the quantity. Thus a Bayesian analysis can be viewed as a generalization of the classical method, with the flexibility to include prior information in a formal way if desired.

Instead of producing confidence intervals, Bayesian analyses produce *credible intervals* (sometimes called probability intervals, and not to be confused with the credibility interval described in Chapter 43). A 95% credible interval from a Bayesian analysis is a summary of the posterior distribution, such that the probability is equal to 95% that the true quantity is within the interval. This is a particularly intuitive way to express uncertainty, and is one of the most appealing aspects of a Bayesian analysis.

We can make further statements about the unknown quantity after a Bayesian analysis. In particular, we can state the probability that the quantity lies in any specified range. For instance, we can state the probability that the quantity is smaller than (or bigger than) zero. This is a bit like a *p*-value, but is a more direct statement since it does not require a null hypothesis for its interpretation.

The controversial aspect of Bayesian meta-analysis is the source of the prior distribution. In practice, several different prior distributions are often compared. If very different prior distributions all lead to the same posterior distribution, then we can conclude that the data are sufficiently convincing to overwhelm any *a priori* belief, and the analysis might be considered robust. However, if the prior is influential, then this usually means that there are insufficient data, and both a Bayesian and a classical meta-analysis ought to be interpreted with some caution.

The combination of prior distribution and data to produce the posterior distribution is computationally very demanding. This is partly why Bayesian methods have only become prominent in recent years. Flexible software is now available for performing Bayesian analyses, particularly the WinBUGS software. In fact, this software is so flexible that very complicated meta-analysis models can be fitted even more easily than in a classical framework. We therefore expect to see more meta-analyses undertaken using Bayesian methods, and for the models they implement to become more and more complex.

### SUMMARY POINTS

- The basic inverse-variance approach to meta-analysis can be applied to a very large class of problems. All we need is a point estimate and its variance from each study.
- Meta-analyses that use the raw data from every included study are often called individual patient data meta-analyses, at least in the medical area. These offer unrivalled flexibility in methods.
- Bayesian methods are based on different philosophical approach to statistics. Results of Bayesian analyses have an appealing interpretation, and they readily allow more complex models to be fitted to meta-analysis datasets, but they require specialized software and require specification of a prior distribution for each parameter.

### Further Reading

- Stewart, L.A. & Tierney, J.F. (2002). To IPD or not to IPD? Advantages and disadvantages of systematic reviews using individual patient data. *Evaluation and the Health Professions* 25: 76–97.
- Sutton, A.J. & Abrams, K.R. (2001). Bayesian methods in meta-analysis and evidence synthesis. *Statistical Methods in Medical Research* 10: 277–303.



## Further Methods



# Overview

Thus far, we have concentrated on a specific approach to meta-analysis, based on obtaining an estimate of effect size, along with its variance, from each study. These effect estimates are combined across studies using a weighted average, with inverse variances as weights. In this Part we consider some alternate approaches.

In Chapter 41 we discuss methods based on combining *p*-values, rather than effect sizes, and also a method based on direction of effect in each study (ignoring the magnitude of the effect). We explain how these methods work, outline their limitations, and discuss when it may be appropriate to apply them.

In Chapter 42 we discuss further methods for dichotomous (binary) outcome data. These include the *Mantel–Haenszel* methods and a ‘one-step’ method known widely as the *Peto* method.

Finally, in Chapter 43 we discuss further methods for correlational data, detailing in particular a method due to Hunter and Schmidt known to many as *psychometric* meta-analysis.



# Meta-Analysis Methods Based on Direction and *p*-Values

---

Introduction  
Vote counting  
The sign test  
Combining *p*-values

---

## INTRODUCTION

In this volume we have concentrated almost exclusively on meta-analysis of *effect sizes*. This reflects current practice, as the overwhelming majority of meta-analyses published over the past two decades *are* meta-analyses of effect sizes. A meta-analysis of effect sizes addresses the magnitude of the effect, which is what we care about. By contrast, a meta-analysis of *p*-values tells us only that the effect is probably not zero.

Nevertheless, there are cases where a meta-analysis of effect sizes is not possible, and in those cases the only option may be to test a null hypothesis. In this chapter we outline a few options for this goal and explain where these might be used.

## VOTE COUNTING

Vote counting is the process of counting the number of studies that are statistically significant and the number that are not, and then choosing the winner. This approach is discussed in Chapter 33 where we explain why it has no validity whatsoever. We mention it here only to be sure that it is not confused with the sign test (below) which is a valid approach.

## THE SIGN TEST

In a sign test, we count the number of studies with findings in one direction and compare this with the number of studies with findings in the other direction, irrespective

of whether the findings were statistically significant. The sign test takes into account *neither* the actual effect magnitudes observed in the studies *nor* the amount of evidence within each study (for example, the sample sizes). As such it has very limited value. However, it might be considered in any of the following cases.

- When no numerical data are provided from the studies, but directions of effect are provided
- When the numerical data are of such different types that they cannot be combined statistically
- When the studies are so diverse in their populations or other characteristics that a pooled effect size is meaningless, but the studies are still addressing a question sufficiently similar that the direction of the effect is meaningful.

If a treatment is truly ineffective, we would expect half of the studies to lie on each side of the no-effect line. We can test this formally by comparing the number of studies in one direction versus the null value of 50%.

For example, in the streptokinase meta-analysis (Chapter 2), 25 out of 33 studies favored the treatment (i.e. had a point estimate of the risk ratio less than 1.0), and 8 studies favored the control (i.e. had a point estimate of the risk ratio greater than 1.0). The two-sided *p*-value for the sign test is 0.00455 (in Excel, the function =2\*BINOMDIST(8,33,0.5,TRUE) returns 0.00455). Or, the one-sided *p*-value for the sign test is 0.0023 (the function =BINOMDIST(8,33,0.5,TRUE) returns 0.0023). Note that in Excel we need to enter the *smaller* of the two numbers and the total (here, 8 and 33).

## COMBINING *p*-VALUES

Another option when we don't want to work with effect sizes is to work directly with the *p*-values from each test, to yield an overall *p*-value. Unlike the sign test or the *p*-value for a summary effect size, both of which test the null hypothesis that the *mean effect* across studies is zero (or 1.0 for a ratio), the tests based on combining *p*-values tests the null hypothesis that the effect size is zero in *all studies*. In other words, if we combine *p*-values and obtain a significant effect, we would conclude that the effect is real *in at least one* of the included studies.

In deciding which approach to use (a meta-analysis of effect sizes or of *p*-values), the fact that some (or all) studies reported *p*-values (and not an effect size) should not be a factor. This is because starting with the *p*-value and some additional information (such as the sample size) we can usually back-compute the effect size and its variance, and then perform a meta-analysis of the effect sizes.

Rather, the approach of combining *p*-values may be considered under the following conditions.

- If we want to test the null hypothesis that the effect is zero *in all the studies*. This might be the case, for example, if each study looked for a different serious side effect for a drug, and we want to know if there is evidence that *any* of the side effects is present. We assume here that the separate tests have adequate power.

- When we have the  $p$ -values but not the sample sizes from each study (and therefore cannot back-compute the effect size)
- When the studies are so diverse in their populations or other characteristics that a pooled effect size is meaningless, but it is meaningful to ask if any of the effects is nonzero.

The last two of these points are the same as those we listed for the sign test, but here we have more information from each study (the  $p$ -value as well as the direction of effect). Since the null hypothesis is different for the sign test (a nonzero effect on average) than for the test of combined  $p$ -values (a nonzero effect in at least one study), we would select the test that matches our null hypothesis.

We describe two methods for performing meta-analyses of  $p$ -values. For both methods it is critical that the starting point for the analysis is a set of exact one-tailed (or one-sided)  $p$ -values. This means that an effect in one direction yields a  $p$ -value in the range of 0.0 to < 0.5, while an effect in the other direction yields a  $p$ -value in the range of > 0.5 to 1.0. A study where the effect was identical in both groups would have a  $p$ -value of exactly 0.50.

If we start with a two-tailed  $p$ -value, we convert this to a one-tailed  $p$ -value as follows. If the effect is in the expected direction, then

$$p_1 = \frac{p_2}{2}. \quad (41.1)$$

If the effect is in the other direction, then

$$p_1 = 1 - \left( \frac{p_2}{2} \right). \quad (41.2)$$

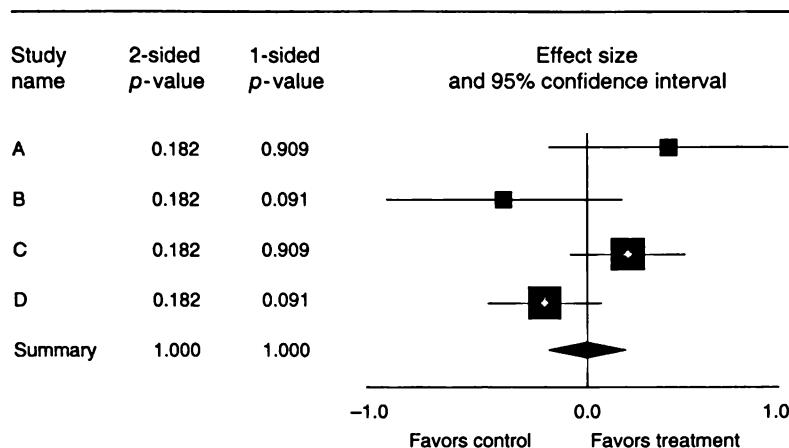
One-tailed  $p$ -values contain information about the direction of effect, whereas two-tailed  $p$ -values do not. For instance, consider the studies in Figure 41.1. The first two studies ( $A$  and  $B$ ) have the same two-sided  $p$ -value, but have effects in opposite directions: the first favors control but the second favors treatment. The one-tailed  $p$ -value reflects this.

The same is true for studies  $C$  and  $D$ . We include these extra studies in the schematic as a reminder that the methods we describe for combining  $p$ -values treat studies  $A$  and  $C$  in exactly the same way, and treat studies  $B$  and  $D$  in exactly the same way. In other words, they do not distinguish between a  $p$ -value arising from a large effect in a small study and the same  $p$ -value arising from a smaller effect in a larger study.

As was true when we were working with effect sizes (and using a  $p$ -value to compute the effect size) we need to start with the actual  $p$ -value. If we are told only that the  $p$ -value falls under 0.05 (for example) we may elect to work with 0.05 or to omit the study from the analysis. Using 0.05 when the actual value could be much lower than 0.05 will lower the chances of a type I error if the null hypothesis is true, but increase the chances of a type II error if the null hypothesis is false.

The first test based on  $p$ -values is known as Fisher's method. We calculate

$$X^2 = -2 \sum_{i=1}^k \ln(p_i), \quad (41.3)$$



**Figure 41.1** Effect size in four fictional studies.

that is, minus 2 times the sum of the logged *p*-values, where  $p_i$  is the one-sided *p*-value from study (*i*) and  $k$  is the number of studies. Under the null hypothesis of no effect in every study,  $X^2$  will follow a central chi-squared distribution with degrees of freedom equal to  $2 \times k$ , so we can report a *p*-value for the aggregated evidence across studies.

The second method is known as Stouffer's method. We calculate a standard normal deviate,  $Z_i$ , from each one-sided *p*-value (this standard normal deviate is often calculated directly from the effect size and its standard error rather than via the *p*-value), and calculate

$$Z_{\text{Stouffer}} = \frac{\sum_{i=1}^k Z_i}{\sqrt{k}}, \quad (41.4)$$

where  $k$  is again the number of studies. Under the null hypothesis of no effect in every study,  $Z_{\text{Stouffer}}$  will follow a standard normal distribution, so we can report a *p*-value for the aggregated evidence across studies.

We will again use the streptokinase example to illustrate these methods (but note that we are testing a different null hypothesis than we had earlier). Table 41.1 presents risk ratios and their 95% confidence intervals for the 33 studies. The Z statistics are obtained as

$$Z_i = \frac{\ln(RR_i)}{(\ln(\text{Upper limit}_i) - \ln(\text{Lower limit}_i))/(2 \times 1.96)}, \quad (41.5)$$

or

$$Z_i = \frac{\ln(RR_i)}{SE_{\ln(RR_i)}}, \quad (41.6)$$

where the numerator is the log risk ratio and the denominator is the standard error of the log risk ratio. The one-tailed *p*-values are obtained by comparing the Z values with

**Table 41.1** Streptokinase data – calculations for meta-analyses of *p*-values.

Study	Risk ratio	Lower limit	Upper limit	Z	Two-tailed <i>p</i> -value	One-tailed <i>p</i> -value	Log(one-tailed <i>p</i> -value)
Fletcher	0.229	0.030	1.750	-1.420	0.155	0.078	-2.554
Dewar	0.571	0.196	1.665	-1.026	0.305	0.152	-1.881
European 1	1.349	0.743	2.451	0.984	0.325	0.837	-0.177
European 2	0.703	0.534	0.925	-2.519	0.012	0.006	-5.134
Heikinheimo	1.223	0.669	2.237	0.654	0.513	0.743	-0.296
Italian	1.011	0.551	1.853	0.034	0.973	0.513	-0.667
Australian 1	0.779	0.478	1.268	-1.005	0.315	0.157	-1.849
Frankfurt 2	0.457	0.252	0.828	-2.581	0.010	0.005	-5.315
NHLBI SMIT	2.377	0.649	8.709	1.307	0.191	0.904	-0.100
Frank	0.964	0.332	2.801	-0.068	0.946	0.473	-0.749
Valere	1.048	0.481	2.282	0.117	0.907	0.547	-0.604
Klein	2.571	0.339	19.481	0.914	0.361	0.820	-0.199
UK-Collab	0.922	0.609	1.394	-0.386	0.699	0.350	-1.051
Austrian	0.608	0.417	0.886	-2.590	0.010	0.005	-5.338
Australian 2	0.702	0.443	1.110	-1.514	0.130	0.065	-2.734
Lasierra	0.282	0.034	2.340	-1.172	0.241	0.121	-2.116
N Ger Collab	1.161	0.840	1.604	0.905	0.366	0.817	-0.202
Witchitz	0.813	0.263	2.506	-0.361	0.718	0.359	-1.025
European 3	0.612	0.356	1.050	-1.782	0.075	0.037	-3.286
ISAM	0.880	0.619	1.250	-0.713	0.476	0.238	-1.436
GISSI-1	0.827	0.749	0.914	-3.738	0.000	0.000	-9.286
Olson	0.429	0.041	4.439	-0.710	0.477	0.239	-1.432
Baroffio	0.079	0.005	1.350	-1.752	0.080	0.040	-3.222
Schreiber	0.333	0.038	2.925	-0.991	0.322	0.161	-1.828
Cribier	1.095	0.073	16.427	0.066	0.948	0.526	-0.642
Sainsous	0.500	0.132	1.887	-1.023	0.306	0.153	-1.876
Durand	0.621	0.151	2.555	-0.660	0.510	0.255	-1.367
White	0.174	0.040	0.761	-2.323	0.020	0.010	-4.596
Bassand	0.604	0.188	1.944	-0.845	0.398	0.199	-1.614
Vlay	0.462	0.048	4.461	-0.668	0.504	0.252	-1.378
Kennedy	0.654	0.322	1.331	-1.171	0.241	0.121	-2.114
ISIS-2	0.769	0.704	0.839	-5.869	0.000	0.000	-19.940
Wisenberg	0.244	0.051	1.164	-1.770	0.077	0.038	-3.260
Sum				-33.677			-89.268

a standard normal distribution. In Excel, the function =NORMSDIST(*Z*) can be used. For example, =NORMSDIST(-1.420) returns 0.078 for the Fletcher study. The sum of the *Z* values is -33.677, and the sum of the logs of the one sided *p*-values is -89.268.

For Fisher's method, we compute

$$\chi^2 = -2 \times -89.268 = 178.5360.$$

The overall *p*-value across studies is obtained by comparing this with a chi-squared distribution with 66 degrees of freedom. In Excel, the function =CHIDIST(178.54,66) returns a very small *p*-value (*p* =  $2.6 \times 10^{-12}$ ).

For Stouffer's method, we compute

$$Z_{\text{Stouffer}} = \frac{-33.677}{\sqrt{33}} = -5.862.$$

The overall  $p$ -value across studies is obtained by comparing this with a standard normal distribution. In Excel, the function =NORMSDIST(-5.862) returns a very small  $p$ -value ( $p = 2.3 \times 10^{-9}$ ).

Both methods provide strong evidence of the benefit of streptokinase in at least one study.

### SUMMARY POINTS

- A meta-analysis of effect sizes is generally the preferred approach since it addresses the issue of interest (*What is the magnitude of the effect?*) rather than the null hypothesis (*Is the effect size zero?*). However, when this approach is not possible, we may consider an approach that tests the null hypothesis.
- The sign test addresses the null hypothesis that the mean effect across all studies is zero. It can be used when we know the direction (but not the magnitude) of the effects, or when the studies are so different that it does not make sense to combine effect sizes.
- Tests to combine  $p$ -values address the null hypothesis that the effect in all studies is zero.

# Further Methods for Dichotomous Data

---

### Introduction

Mantel–Haenszel method

One-step (Peto) formula for odds ratio

---

## INTRODUCTION

In this chapter we present two methods, the Mantel–Haenszel method and the one-step method (also known as the Peto method) for performing a meta-analysis on odds ratios. For both methods we assume the data from each study are presented in the form of a  $2 \times 2$  table with cells labeled as in Table 42.1.

## MANTEL–HAENSZEL METHOD

Widely familiar to epidemiologists, although perhaps less familiar to others, the Mantel–Haenszel (*MH*) method is unusual in being a weighted average of odds ratios rather than of log odds ratios. If we use  $Y_i$  to denote the odds ratio in study  $i$ , then  $Y_i$  is computed as

$$Y_i = \frac{A_i D_i}{B_i C_i}. \quad (42.1)$$

In the Mantel–Haenszel method, the weight assigned to each study is

$$W_i = \frac{B_i C_i}{n_i}, \quad (42.2)$$

where

$$n_i = A_i + B_i + C_i + D_i, \quad (42.3)$$

**Table 42.1** Nomenclature for  $2 \times 2$  table of events by treatment.

	Events	Non-events	
Treated	A	B	$n_1$
Control	C	D	$n_2$

and the weighted mean is then computed as

$$OR_{MH} = \frac{\sum_{i=1}^k W_i Y_i}{\sum_{i=1}^k W_i}. \quad (42.4)$$

This is the sum of the products (effect size multiplied by weight) divided by the sum of the weights. This formula is identical to the one for inverse-variance, but the weights ( $W$ ) have been defined differently and the effect size ( $Y$ ) is in raw units rather than log units.

While the odds ratio itself is computed in raw units, the variance is computed in log units. Therefore we need to transform the odds ratio into log units to compute the Z-score and confidence intervals. The natural log of the  $MH$  odds ratio is simply

$$\ln OR_{MH} = \ln(OR_{MH}). \quad (42.5)$$

Recall that in the inverse variance formula the variance of the summary effect was defined as the reciprocal of the sum of the weights. This is not the case here. Rather, the  $MH$  approach calls for us to accumulate separate values, which will be summed across studies and then used to compute the variance of the summary effect.

For each study ( $i$ ),

$$R_i = \frac{A_i D_i}{n_i}, \quad (42.6)$$

$$S_i = \frac{B_i C_i}{n_i}, \quad (42.7)$$

$$E_i = \frac{(A_i + D_i) A_i D_i}{n_i^2}, \quad (42.8)$$

$$F_i = \frac{(A_i + D_i) B_i C_i}{n_i^2}, \quad (42.9)$$

$$G_i = \frac{(B_i + C_i) A_i D_i}{n_i^2}, \quad (42.10)$$

and

$$H_i = \frac{(B_i + C_i) B_i C_i}{n_i^2}. \quad (42.11)$$

Then the variance of the summary effect, in log units, is

$$V_{\ln OR_{MH}} = 0.5 \left( \frac{\sum_{i=1}^k E_i}{\left( \sum_{i=1}^k R_i \right)^2} + \frac{\sum_{i=1}^k F_i + \sum_{i=1}^k G_i}{\sum_{i=1}^k R_i \times \sum_{i=1}^k S_i} + \frac{\sum_{i=1}^k H_i}{\left( \sum_{i=1}^k S_i \right)^2} \right) \quad (42.12)$$

and

$$SE_{\ln OR_{MH}} = \sqrt{V_{\ln OR_{MH}}}. \quad (42.13)$$

The 95% confidence interval for the summary effect in log units would be computed as

$$LL_{\ln OR_{MH}} = \ln OR_{MH} - 1.96 \times SE_{\ln OR_{MH}} \quad (42.14)$$

and

$$UL_{\ln OR_{MH}} = \ln OR_{MH} + 1.96 \times SE_{\ln OR_{MH}}. \quad (42.15)$$

The Z-value is given by

$$Z = \frac{\ln OR_{MH}}{SE_{\ln OR_{MH}}}. \quad (42.16)$$

For a one-tailed test the p-value is given by

$$p = 1 - \Phi(\pm|Z|), \quad (42.17)$$

where we choose '+' if the difference is in the expected direction and '-' otherwise. Or, for a two-tailed test by

$$p = 2[1 - (\Phi(\pm|Z|))]. \quad (42.18)$$

We present a worked example for Dataset 2, which was originally presented in Chapter 18 and served as an example for the inverse-variance method.

From the results in Table 42.2, the Mantel–Haenszel summary odds ratio is computed as

$$OR_{MH} = \frac{31.5452}{67.0452} = 0.4705,$$

**Table 42.2** Mantel–Haenszel – odds ratio.

	Cell counts					Compute odds ratio		
	A	B	C	D	N	Odds ratio	Weight	OR*WT
Saint	12	53	16	49	130	0.6934	6.5231	4.5231
Kelly	8	32	10	30	80	0.7500	4.0000	3.0000
Pilbeam	14	66	19	61	160	0.6810	7.8375	5.3375
Lane	25	375	80	320	800	0.2667	37.5000	10.0000
Wright	8	32	11	29	80	0.6591	4.4000	2.9000
Day	16	49	18	47	130	0.8526	6.7846	5.7846
Sum							67.0452	31.5452

and the log odds ratio is

$$\ln OR_{MH} = \ln (0.4705) = -0.7539.$$

The confidence interval for the summary effect is calculated in log units as follows (see Table 42.3).

For the first study (Saint), for example,

$$R = \frac{(12)(49)}{130} = 4.5231,$$

$$S = \frac{(53)(16)}{130} = 6.5231,$$

$$E = \frac{(12 + 49)(12)(49)}{130^2} = 2.1224,$$

$$F = \frac{(12 + 49)(53)(16)}{130^2} = 3.0608,$$

$$G = \frac{(53 + 16)(12)(49)}{130^2} = 2.4007,$$

and

$$H = \frac{(53 + 16)(53)(16)}{130^2} = 3.4622.$$

The variance of the log odds ratio is given by

$$V_{\ln OR_{MH}} = 0.5 \times \left( \frac{14.5064}{31.5452^2} + \frac{30.1295 + 17.0388}{31.5452 \times 67.0452} + \frac{36.9157}{67.0452^2} \right) = 0.0225,$$

so

$$SE_{\ln OR_{MH}} = \sqrt{0.0225} = 0.1502,$$

$$LL_{\ln OR_{MH}} = -0.7539 - 1.96 \times 0.1502 = -1.0482,$$

$$UL_{\ln OR_{MH}} = -0.7539 + 1.96 \times 0.1502 = -0.45951,$$

and

$$Z = \frac{-0.7539}{0.1502} = -5.0211.$$

**Table 42.3** Mantel–Haenszel – variance of summary effect.

	Cell counts					Compute variance of combined effect					
	A	B	C	D	N	R	S	E	F	G	H
Saint	12	53	16	49	130	4.5231	6.5231	2.1224	3.0608	2.4007	3.4622
Kelly	8	32	10	30	80	3.0000	4.0000	1.4250	1.9000	1.5750	2.1000
Pilbeam	14	66	19	61	160	5.3375	7.8375	2.5020	3.6738	2.8355	4.1637
Lane	25	375	80	320	800	10.0000	37.5000	4.3125	16.1719	5.6875	21.3281
Wright	8	32	11	29	80	2.9000	4.4000	1.3413	2.0350	1.5588	2.3650
Day	16	49	18	47	130	5.7846	6.7846	2.8033	3.2879	2.9813	3.4967
Sum						31.5452	67.0452	14.5064	30.1295	17.0388	36.9157

For a one-tailed test  $p$  is given by

$$p = 1 - \Phi(+5.0211) = 0.00000026,$$

which, in Excel is  $=1-(\text{NORMSDIST}(\text{ABS}(Z)))$ . Or, for a two-tailed test,  $p$  is given by

$$p = 2[1 - (\Phi(|-5.0211|))] = 0.00000051$$

which, in Excel is  $=1-(\text{NORMSDIST}(\text{ABS}(Z)))*2$ .

Finally, we would convert the values back to raw units for display, using

$$OR_{MH} = 0.4705,$$

$$LL_{OR_{MH}} = \exp(-1.0482) = 0.3506,$$

and

$$UL_{OR_{MH}} = \exp(0.4596) = 0.6315.$$

The  $Z$ -value and  $p$ -value are the same as for the log values.

The Mantel–Haenszel method was developed for combining odds ratios across  $2 \times 2$  tables. It has since been extended by others to combine risk ratios or risk differences across  $2 \times 2$  tables. Similar formulas are available for these, although we do not present them since they provide little insight into the pooling mechanism and they can readily be implemented in meta-analysis software (see Chapters 49 and 51).

The Mantel–Haenszel method is based on the fixed-effect model, where the weight assigned to each study is based on that study alone and not on the variance across studies.

### ONE-STEP (PETO) FORMULA FOR ODDS RATIO

The one-step method for computing the summary odds ratio works on the log odds ratio scale, and is a variant of the basic inverse-variance approach. However, this method uses a different formula than the one presented earlier for both the odds ratio and its variance. The one-step method is sometimes called the Peto method.

The log odds ratio in study ( $i$ ) is estimated using

$$Y_i = \frac{O_i - E_i}{I_i}, \quad (42.19)$$

where the observed count,  $O_i$  is given by

$$O_i = A_i, \quad (42.20)$$

and the expected count,  $E_i$ , is given by

$$E_i = \frac{(A_i + B_i) \times (A_i + C_i)}{n_i}, \quad (42.21)$$

where

$$n_i = A_i + B_i + C_i + D_i, \quad (42.22)$$

and

$$I_i = \frac{(A_i + B_i) \times (C_i + D_i) \times (A_i + C_i) \times (B_i + D_i)}{n_i^2 \times (n_i - 1)}. \quad (42.23)$$

The variance of the log odds ratio estimate for a single study is

$$V_{Y_i} = \frac{1}{I_i}, \quad (42.24)$$

where the weight given to the study is its inverse-variance given by

$$W_i = \frac{1}{V_{Y_i}}. \quad (42.25)$$

Note that the weight is just  $I_i$ . The meta-analysis is given by the weighted average of the log odds ratio estimates

$$\ln OR_{onestep} = \frac{\sum_{i=1}^k W_i Y_i}{\sum_{i=1}^k W_i} \quad (42.26)$$

with variance given by

$$V_{\ln OR_{onestep}} = \frac{1}{\sum_{i=1}^k W_i}. \quad (42.27)$$

An alternative way of writing these results is

$$\ln OR_{onestep} = \frac{\sum_{i=1}^k (O_i - E_i)}{\sum_{i=1}^k I_i}, \quad (42.28)$$

with variance given by

$$V_{\ln OR_{onestep}} = \frac{1}{\sum_{i=1}^k I_i}. \quad (42.29)$$

The standard error (in log units) is

$$SE_{\ln OR_{onestep}} = \sqrt{V_{\ln OR_{onestep}}}. \quad (42.30)$$

The 95% confidence interval for the summary effect in log units would be computed as follows.

$$LL_{\ln OR_{onestep}} = \ln OR_{onestep} - 1.96 \times SE_{\ln OR_{onestep}} \quad (42.31)$$

and

$$UL_{\ln OR_{onestep}} = \ln OR_{onestep} + 1.96 \times SE_{\ln OR_{onestep}}. \quad (42.32)$$

The Z-value is given by

$$Z = \frac{\ln OR_{onestep}}{SE_{\ln OR_{onestep}}}. \quad (42.33)$$

The p-value for a one-tailed test given by

$$p = 1 - \Phi(\pm|Z|), \quad (42.34)$$

where we choose '+' if the difference is in the expected direction and '-' otherwise. Or, for a two-tailed test as

$$p = 2[1 - (\Phi(|Z|))]. \quad (42.35)$$

In Excel, the two-tailed p-value is given by  $= (1 - (NORMSDIST(ABS(Z)))) * 2$ .

In Table 42.4 we apply these formulas to Dataset 2 (see Table 14.4). Then, using sums computed in Table 42.4, we compute

$$\ln OR_{onestep} = \frac{-35.5000}{48.4827} = -0.7322,$$

$$V_{\ln OR_{onestep}} = \frac{1}{48.4827} = 0.0206,$$

$$SE_{\ln OR_{onestep}} = \sqrt{0.0206} = 0.1436,$$

$$LL_{\ln OR_{onestep}} = -0.7322 - 1.96 \times 0.1436 = -1.0137,$$

$$UL_{\ln OR_{onestep}} = -0.7322 + 1.96 \times 0.1436 = -0.4507,$$

$$Z = \frac{-0.7322}{0.1436} = -5.098,$$

$$p = 1 - \Phi(-5.098) = 0.00000017,$$

and

$$p = 2[1 - (\Phi(|-5.098|))] = 0.00000034,$$

Finally, we would convert the values back to raw units for display, using

$$OR_{onestep} = \exp(-0.7322) = 0.481,$$

$$LL_{OR_{onestep}} = \exp(-1.0137) = 0.363,$$

**Table 42.4** One-step – odds ratio and variance.

	A	B	C	D	N	O	E	O-E	I	Y	W	WY
Saint	12	53	16	49	130	12	14	-2	5.5349	-0.3613	5.5349	-2.0000
Kelly	8	32	10	30	80	8	9	-1	3.5316	-0.2832	3.5316	-1.0000
Pilbeam	14	66	19	61	160	14	16.5	-2.5	6.5896	-0.3794	6.5896	-2.5000
Lane	25	375	80	320	800	25	52.5	-27.5	22.8332	-1.2044	22.8332	-27.5000
Wright	8	32	11	29	80	8	9.5	-1.5	3.6677	-0.4090	3.6677	-1.5000
Day	16	49	18	47	130	16	17	-1	6.3256	-0.1581	6.3256	-1.0000
Sum											48.4827	35.5000

and

$$UL_{OR_{onestep}} = \exp(-0.4507) = 0.637,$$

In the one-step approach, like the basic approach, all analyses are carried out on the log of the odds ratio. We compute a weighted mean of the log values, and then exponentiate this value to report the summary effect.

When one or more of the cells in the  $2 \times 2$  table is empty (that is, a value of zero), the basic inverse-variance formula cannot work with zero, and the typical approach is to add the value 0.5 (or some other value) to all four cells. However, both the Mantel-Haenszel method and the one-step approach are able to work with a value of zero, and so no adjustment is needed.

The one-step approach follows the same logic as the basic scheme, in that it uses inverse-variance weights at all points in the analysis, and differs from the *basic* approach only in that it uses a slightly different way of computing the log odds ratio and its variance. As such, the one-step approach can be extended to the random-effects model. In practice, though, this is rarely done.

### SUMMARY POINTS

- The Mantel-Haenszel method for combining odds ratios is an alternative to the fixed-effect inverse variance method.
- The one-step (Peto) method for combining odds ratios is an inverse-variance method, but uses an alternate approach to computing the odds ratio and variance in each study. This method offers some advantages when some studies have empty cells.

# Psychometric Meta-Analysis

---

### Introduction

The attenuating effects of artifacts

### Meta-analysis methods

Example of psychometric meta-analysis

Comparison of artifact correction with meta-regression

Sources of information about artifact values

How heterogeneity is assessed

Reporting in psychometric meta-analysis

Concluding remarks

---

## INTRODUCTION

Most meta-analyses aim to summarize the results obtained in studies that were carried out with an implicit assumption that those results are the *best available* estimates of effect. However, any study has methodological flaws that affect its results. If we were able to *correct* estimates for these flaws, it would be preferable to perform a meta-analysis of these corrected results. Then we can address what the study results *would have been* if all of the studies had been free of methodological imperfections (including finite sample size) and to estimate parameters describing the effects in these methodologically perfect studies.

Unfortunately it is typically difficult, if not impossible, to know the specific impact of these methodological flaws. A growing literature addresses attempts to determine the likely biases using theoretical considerations or information from other studies and other meta-analyses. One approach to meta-analysis has focused since its inception almost exclusively on developing ways to adjust estimates of effect for methodological limitations of the studies. This is the field of psychometric meta-analysis (also called validity generalization or Hunter–Schmidt meta-analysis). This field has also adopted a somewhat different (though closely related) set of methods for combining results in meta-analysis than those that have been discussed previously. This chapter provides an overview of two issues. One is the approach to estimates of effect (known

as artifact correction), which will be of interest to nearly anyone thinking about using meta-analysis. The other is the methods that are commonly used to combine results in the field of psychometric meta-analysis, which will be of interest primarily to researchers who use correlations as their effect size measure.

### **Psychometric meta-analysis**

Methodological flaws in research studies affect study results in ways that might be thought of as artifacts of the study design. Much of the work on psychometric meta-analysis has focused on the case of studies that use continuous outcome measures and use effect size measures involving standardization (such as correlation coefficients or standardized mean differences). Since the measuring instruments are subject to measurement error (imperfect reliability), they produce effect size estimates that are made smaller (attenuated) by this measurement error. If the samples in some studies are selected in ways that do not contain the full range of variation on either independent or dependent variables, effect size estimates are reduced (or attenuated) due to this restriction of range. The dichotomization of variables that are inherently continuous in order to form binary categories produces effect size estimates that are attenuated by the reduction in variance that results from the (artificial) dichotomization. A more complete list of methodological problems that can have an impact on observed effects may be found in Hunter and Schmidt (2004). Psychometric meta-analysts rightly argue that what researchers really would like to know about is the relationship between the constructs (or variables) in a study that are *not* artifacts of study design. Therefore the answers to scientific questions are best provided by estimating the results that would have been observed had each study been free of methodological imperfections. This is what psychometric meta-analysis would call the *true values*.

Psychometric meta-analysis uses procedures based on psychometric principles (hence the name psychometric meta-analysis) to make corrections for attenuation due to measurement error and other artifacts at the level of the individual effect size estimate, before these effects are synthesized across the set of studies in the meta-analysis.

Many methodological artifacts have the effect of attenuating relationships among variables. Therefore, the average effect size corrected for the effects of artifacts (such as measurement unreliability or range restriction) will generally be larger than if the corrections were not made.

Our estimate of the *variation* in effects may also be different. The reason is that the impact of artifacts on effect sizes varies from study to study, and this increases the variability of the observed effect sizes. Consequently, the variation across studies in the effects corrected for artifacts (the variation that is of interest in psychometric meta-analyses) is typically less than that in the observed effects.

### **THE ATTENUATING EFFECTS OF ARTIFACTS**

Psychometric meta-analyses usually represent the relationship between the observed (unadjusted) effect size and the *true* (adjusted) effect size via their ratio. In accordance

with psychometric terminology, however, they use the term *attenuated* to refer to the observed effect, and the term *unattenuated* to refer to the adjusted effect. The ratio of the attenuated to the *unattenuated* effect describes the impact of the artifact on the effect size, and is called an *artifact multiplier* because the magnitude of the observed (attenuated) effect size is equal to the artifact multiplier times the *unattenuated* effect size. For example, if the effect size is a correlation coefficient (as it typically is in psychometric meta-analyses), and the artifact is measurement error (unreliability) in one of the two variables, then the ratio  $a$  of the attenuated correlation  $\rho$  to the unattenuated correlation  $\rho^u$

$$a = \frac{\rho}{\rho^u} \quad (43.1)$$

is the square root of the reliability coefficient of  $Y$ . (The reliability coefficient is an index used in psychometric theory and other areas to characterize the reproducibility of measurements.) Although a discussion of measurement theory in general, and reliability, in particular, are beyond the scope of this book, a sophisticated discussion of reliability can be found in Lord and Novick (1968) and an introductory discussion can be found in Crocker and Algina (1986). Other artifacts (such as restriction of the range of measurement for one or other variable) lead to different expressions for the artifact multiplier. Algebraic expressions for the artifact multipliers come from psychometric or statistical theory (see Hunter and Schmidt, 2004). While artifacts typically lead to attenuation (that is,  $0 < a < 1$ ), this need not be so for all artifacts (e.g. in the case where the range of variables in the study population is greater than that in the reference population of interest).

If several artifacts influence an effect size parameter, we can compute an artifact multiplier for each. Then the combined effect of all of the artifacts can be expressed by a combined artifact multiplier that is the product of all of the individual artifact multipliers. This combined artifact multiplier can be used like any single artifact multiplier. The object of psychometric meta-analysis is to describe the distribution of the unattenuated effect size parameters (such as the  $\rho^u$ ) given estimates of the observed (attenuated) effect size estimates (such as the  $r$ ). It follows from (43.1) that the unattenuated correlation can be estimated from the observed correlation and the artifact multiplier as

$$r^u = \frac{r}{a}. \quad (43.2)$$

Because the artifact multipliers are taken to be constants, it follows that the variance of  $r^u$  (call this  $V_r^u$ ) is  $1/a^2$  as large as the variance of  $r$  (call this  $V_r$ ). That is,

$$V_r^u = \frac{V_r}{a^2}. \quad (43.3)$$

Having obtained estimates of unattenuated correlations and their variances, we could implement any of the methods described in Parts 3 to 6. For instance, we could perform a fixed-effect or a random-effects meta-analysis using basic inverse-variance weighted averages, we could assess heterogeneity using a statistical test, or  $I^2$ , or by computing  $T^2$ , and could perform subgroup analyses and meta-regression to explore heterogeneity. The methods that are typically used in the psychometric meta-analysis field are somewhat different, however. In particular,

- Raw correlations are used as the effect size index instead of Fisher's Z-transformed correlations.
- Sample sizes are used as weights instead of inverse variances (which makes little difference for correlations, but which could yield very different results, were it applied to binary data).
- A different method is used to estimate  $\tau^2$ .
- Heterogeneity is examined and reported in a different way.

## META-ANALYSIS METHODS

We will describe first the methods used to perform a meta-analysis of the observed (attenuated) correlations, and then apply similar ideas to the unattenuated correlations. At this point we add subscripts ( $i$ ) to refer to the different studies, and use  $k$  to denote the number of studies. The convention in psychometric meta-analysis is to use the sample sizes as weights for computing the mean correlation. This yields the mean

$$\bar{r} = \frac{\sum_{i=1}^k n_i r_i}{\sum_{i=1}^k n_i}. \quad (43.4)$$

To estimate the between-studies variance ( $\tau^2$ ) of the underlying attenuated correlation coefficients, we require a variance for the correlation from each study. We could use the usual variance estimate in (6.1). However, the convention in psychometric meta-analysis is to use the sample size weighted estimate of the mean effect size parameter to compute the sampling error variance, namely

$$V_{r_i} = \frac{[1 - (\bar{r})^2]^2}{n_i - 1}. \quad (43.5)$$

To compute the between-studies variance component of the observed (attenuated) effect size parameters we first compute the sample-size-weighted variance of the observed (attenuated) correlations

$$S^2 = \frac{\sum_{i=1}^k n_i (r_i - \bar{r})^2}{\sum_{i=1}^k n_i}. \quad (43.6)$$

Then the between-studies variance component of the observed (attenuated) effect size parameters is computed as the difference between  $S^2$  and the sample-size-weighted average of the sampling error variances, namely

$$T^2 = S^2 - \frac{\sum_{i=1}^k n_i V_{r_i}}{\sum_{i=1}^k n_i}. \quad (43.7)$$

Note that the right hand side of this equation consists of two terms. The first term is a weighted variance of the observed correlations (with the  $i^{\text{th}}$  correlation weighted by  $n_i$ ) and the second term is a weighted average of the sampling error variances using the same weights. The between-studies variance estimate is the difference between an observed (weighted) variance and the (weighted) average of the sampling error variances. Thus the between-studies variance estimate can be seen as the observed variance *adjusting for* the sampling error variance.

A meta-analysis of unattenuated (artifact-corrected) correlation estimates follows a similar procedure with  $r_i^u$  replacing  $r_i$  and with revised weights. Recall that the artifact correction for the variance in each study is given by (43.3). The analogous correction to the sample size is to adjust it from  $n_i$  to  $a_i^2 n_i$ .

The mean unattenuated correlation is therefore

$$\bar{r}^u = \frac{\sum_{i=1}^k n_i a_i^2 r_i^u}{\sum_{i=1}^k n_i a_i^2}. \quad (43.8)$$

The weighted variance of the unattenuated correlations is

$$(S^u)^2 = \frac{\sum_{i=1}^k n_i a_i^2 (r_i^u - \bar{r}^u)^2}{\sum_{i=1}^k n_i a_i^2}. \quad (43.9)$$

Then the between-studies variance component of the unattenuated effect size parameters  $(T^u)^2$  is computed as the difference between  $(S^u)^2$  and the sample-size-weighted average of the sampling error variances of the unattenuated correlations, namely

$$(T^u)^2 = (S^u)^2 - \frac{\sum_{i=1}^k n_i a_i^2 V_i^u}{\sum_{i=1}^k n_i a_i^2}. \quad (43.10)$$

## EXAMPLE OF PSYCHOMETRIC META-ANALYSIS

We now present a hypothetical example to illustrate the methods of psychometric meta-analysis. We will correct for only a single artifact, error of measurement (unreliability) in the dependent variable. In actual applications of psychometric meta-analysis, corrections for multiple artifacts (for example, error of measurement in both independent and dependent variables, and restriction of range in one or both variables) might be used.

Suppose that six studies were conducted to assess the validity of a pre-hire work sample test (independent variable) to predict the job performance (dependent variable) of dental hygienists six months after hire. The first two studies were conducted by a

**Table 43.1** Fictional data for psychometric meta-analysis.

Study	<i>n</i>	<i>r</i>	Criterion reliability
University 1	130	0.24	0.75
University 2	90	0.11	0.75
Private 1	30	0.05	0.60
Private 2	25	0.17	0.60
Volunteer 1	50	0.38	0.90
Volunteer 2	65	0.50	0.90

consortium of clinics run by schools of dentistry, and job performance was measured using a standard, professionally developed rating scale. The third and fourth studies were conducted in large multi-partner private dental practices, and job performance was measured using a home-grown rating scale developed by a group of the partners from the two practices. The fifth and sixth studies were done at clinics run by a non-profit organization where dentists volunteered to work for free two weeks per year. Job performance was measured using a standardized work behavior assessment scale designed by the nonprofit organization.

The sample sizes (*n*), the observed (attenuated) correlations (*r*), and the criterion reliabilities are given in Table 43.1. We first perform a meta-analysis of the correlations as they were observed (the attenuated correlations). The calculations are provided in Table 43.2.

We use the first three columns in Table 43.2 to compute the sum of  $n_i$  and the sum of  $n_i r_i$ . Then, the mean correlation is computed as

$$\bar{r} = \frac{\sum_{i=1}^6 n_i r_i}{\sum_{i=1}^6 n_i} = \frac{98.35}{390} = 0.2518,$$

which is inserted into column 4 in the table. Then, we complete the remaining columns in the table.

**Table 43.2** Observed (attenuated) correlations.

Study	<i>n<sub>i</sub></i>	<i>r<sub>i</sub></i>	<i>n<sub>i</sub> r<sub>i</sub></i>	$\bar{r}$	<i>V<sub>ri</sub></i>	$n_i(r_i - \bar{r})^2$	$n_i V_{ri}$
University 1	130	0.24	31.20	0.252179	0.0068	0.019283	0.88365
University 2	90	0.11	9.90	0.252179	0.0099	1.819338	0.88671
Private 1	30	0.05	1.50	0.252179	0.0302	1.226290	0.90709
Private 2	25	0.17	4.25	0.252179	0.0365	0.168835	0.91339
Volunteer 1	50	0.38	19.00	0.252179	0.0179	0.816910	0.89475
Volunteer 2	65	0.50	32.50	0.252179	0.0137	3.991991	0.89056
Total	390		98.35			8.042647	5.37615

The variance of the correlation in study number  $i$  is

$$V_{r_i} = \frac{[1 - 0.25218^2]^2}{n_i - 1}.$$

and these values are listed in Table 43.2.

The between-studies variance is computed as follows:

$$S^2 = \frac{8.0426}{390} = 0.0206$$

and

$$T^2 = 0.0206 - \frac{5.3762}{390} = 0.0069.$$

Calculations for the unattenuated correlations appear in Table 43.3. The artifact multiplier for this particular example is the square root of the criterion reliability.

The mean unattenuated correlation is

$$\bar{r}^u = \frac{88.9048}{301.5} = 0.2949.$$

The between-studies variance is computed as

$$(S^u)^2 = \frac{6.6287}{301.5} = 0.0212$$

and

$$(T^u)^2 = 0.0220 - \frac{5.3762}{301.5} = 0.0042.$$

The effect of the artifact of error of measurement in the criterion variable was to decrease the estimate of the average correlation from about 0.295 to about 0.252 or about 15%. In addition the artifact of error of measurement in the criterion variable increased the estimated variance of the unattenuated correlations compared to the observed (attenuated) correlations from about 0.0042 to about 0.0069 or about 64%. This example illustrates a common finding in psychometric meta-analyses: artifacts tend to reduce the magnitude of effects and increase their apparent variation in the sense that the unattenuated effect parameters are estimated to have a larger mean and smaller between-studies variance than the observed effects.

**Table 43.3** Unattenuated correlations.

Study	$n_i$	$a_i$	$r_i^u$	$n_i a_i^2 r_i^u$	$V_{r_i}^u$	$n_i a_i^2$	$n_i a_i^2 (r_i^u - \bar{r}^u)^2$	$n_i a_i^2 V_{r_i}^u$
University 1	130	0.866	0.28	27.020	0.00906	97.500	0.030707	0.0068
University 2	90	0.866	0.13	8.574	0.01314	67.500	1.901896	0.0099
Private 1	30	0.775	0.06	1.162	0.05039	18.000	0.954894	0.0302
Private 2	25	0.775	0.22	3.292	0.06089	15.000	0.085290	0.0365
Volunteer 1	50	0.949	0.40	18.025	0.01988	45.000	0.502575	0.0179
Volunteer 2	65	0.949	0.53	30.832	0.01522	58.500	3.153359	0.0137
Total	390			88.905		301.500	6.628722	5.3762

### Explained variance in psychometric meta-analyses

It is conventional in psychometric meta-analyses to consider  $S^2$  as the (estimated total) variance of the observed effect estimates and  $(T^u)^2$  as the (estimated) variance of the true effects remaining after the artifacts effects have been removed. Thus

$$\frac{S^2 - (T^u)^2}{S^2} = 1 - \frac{(T^u)^2}{S^2} \quad (43.11)$$

is the proportion of variance in the observed (attenuated) effect estimates explained by artifacts.

Note that this definition is in the same spirit as the index  $I^2$  in conventional meta-analysis, but  $I^2$  focuses on unexplained, as opposed to explained variance. Thus

$$\frac{(T^u)^2}{S^2} \times 100\% \quad (43.12)$$

estimates a quantity similar to  $I^2$  that reflects the percentage of the variance in the observed correlations that is *due to* the variance in the unattenuated correlation parameters.

In our example, the psychometric analysis would estimate the proportion of explained variance as

$$\frac{0.02062 - 0.0042}{0.0206} = 0.7987,$$

so that 79.9% of the total variance in the observed (attenuated) effect size estimates is explained by the artifacts of sampling error and measurement error in the criterion variable.

### COMPARISON OF ARTIFACT CORRECTION WITH META-REGRESSION

In conventional meta-analysis, if criterion unreliability was hypothesized to influence effect size, it would be treated as a covariate. After the core meta-analysis was run, criterion unreliability would be regressed on effect size, to test whether it could explain some of the between-study variation in the (uncorrected) correlations. In this example, we conducted a random-effects meta-analysis on the (uncorrected) dental hygienist validity data. The overall mean effect was computed as 0.25,  $T^2$  was computed 0.012, and  $I^2$  was 44.481. We then performed a random-effects meta-regression (method of moments), regressing criterion unreliability on effect size. Results showed that the slope was significant with a  $p$ -value of 0.011. More to the point,  $T^2$  was 0.00, meaning that all of the explainable variance was accounted for by differences in criterion reliability.

### SOURCES OF INFORMATION ABOUT ARTIFACT VALUES

The most appropriate method of adjusting for the effects of artifacts is to correct each effect individually, using reliability and range restriction information that is provided in the study from which the effect is extracted. In most cases, however, this information is not available. This presents what is essentially a missing data problem. As in other

cases of missing data, a variety of imputation techniques can be used to estimate the missing reliability and/or range restriction values. One technique that is recommended by the developers of psychometric meta-analysis is to create artifact distributions based on the information provided in the subset of studies that report the relevant data. In this case, the meta-analysis is conducted on the uncorrected effects and the average artifact value (e.g. reliability) from the artifact distribution is used to correct the average effect, while the standard deviation of the artifact values is used to correct the variance of the observed effects. The use of hypothetical distributions of artifacts in general, and in particular, the application of distributions created for meta-analyses of employment test validities to quite different research questions, has been criticized by many methodologists who are generally supportive of the psychometric meta-analysis framework. An alternative, when information about measurement reliability and range restriction is largely missing from the set of studies in the analysis, is to use the mean observed effect as the estimate of the population mean effect, and to remove the variance due to sampling error from the total observed variance, to produce an estimate of true variance in effects. Psychometric meta-analysis refers to this alternative as a *bare-bones* meta-analysis.

## HOW HETEROGENEITY IS ASSESSED

### **When covariates are not hypothesized in advance of the meta-analysis**

Psychometric meta-analyses do not usually use the chi-square test of homogeneity to test whether the observed variance is greater than the amount that would be expected due to sampling error. Instead, they have substituted the 75% decision rule, which proposes that if 75% or more of the total (observed) variance is due to artifacts, including sampling error, the researcher may conclude that, actually, all of the variance is artifactual, since there are several commonly operating artifacts for which no corrections can be made (such as transcriptional and coding errors, which are claimed to be ubiquitous). This rule was formulated for the original application of psychometric meta-analysis, the assessment of the consistency of employment test validities, and was tested through computer simulation of conditions typical of employment testing research. Its performance in other research areas remains unstudied, and it is not advisable to use this rule outside the area for which it was developed. Many users of psychometric meta analysis have objected to the 75% rule because it focuses on the percentage of true variance rather than on its magnitude, and may cause substantial remaining true variance to be ignored.

### **When there is an *a priori* hypothesis that a covariate may explain heterogeneity**

In the cases of a hypothesized discrete covariate, the procedure followed by psychometric meta-analysis is to divide the studies into subgroups based on values of the hypothesized covariate. For example, if the hypothesis that the correlations will be higher for males than for females, studies of males will comprise one subgroup, and

studies of females will comprise a second subgroup. Subgroup meta-analyses are conducted, and are declared to be different when the *true* means in each subgroup are different from each other, and the confidence intervals around each mean are largely non-overlapping.

For continuous covariates, the usual meta-regression procedures are followed, using correlations that are individually corrected for artifacts, and the corrected sampling error variances.

### REPORTING IN PSYCHOMETRIC META-ANALYSIS

The findings of a psychometric meta-analysis focus on the value of three parameters: (1) the average percentage of the total (observed) variance of effects across studies that is explained by statistical and measurement artifacts, including sampling error; (2) the estimated mean true effect parameter; and (3) the estimated standard deviation of true effects. In a psychometric meta-analysis this last value represents the degree of dispersion across the true effects, that is, the degree of dispersion remaining after sampling error variance, and variance due to other artifacts, have been removed from the observed variance. It is used to form what psychometric meta-analysis refers to as the credibility interval. The credibility interval contains the distribution of true effects and is roughly analogous to the prediction interval discussed in Chapter 17. Typically, a psychometric meta-analysis results table presents the lower end of the 80 or 90% credibility interval, which is the estimated value of the effect above which 80 or 90% of true effects are expected to be found. The lower bound value of the credibility interval (because we care about the *minimum validity*) and the width of the interval are considered to be the most important results of the meta-analysis.

### CONCLUDING REMARKS

The purpose of this chapter has been to explain the basic tenets of psychometric meta-analysis, and to explain how these differ from other methods of meta-analysis. A key characteristic of psychometric meta-analyses is the correction of individual study results for artifacts, so that we can estimate what the effect would be if there were no methodological limitations. This would be a desirable aim in any meta-analysis. Unfortunately methods are not established for many types of data, and even among supporters of psychometric meta-analysis, there is some disagreement about the specific operational procedures to be followed in making these corrections. Another characteristic of psychometric meta-analyses is the emphasis on credibility intervals rather than overall means. We have argued in Chapter 17 that analogous prediction intervals should be considered routinely for all meta-analyses.

Other differences are technical in nature. These include the way weights are assigned, the way between-studies variance is estimated, the use of correlation coefficients rather than Fisher's Z-transformed values, and the use of the average versus the individual effect in calculating the sampling error variance. The effects of

these choices can be viewed entirely separately from other aspects of the procedures, and have been extensively examined through simulations, which suggest that the differences are likely to be trivial in many cases.

Finally, there are issues that remain the subject of disagreements among users of psychometric meta-analysis. These include the degree to which the assumptions needed to make the corrections are met in specific situations, the use of artifact distributions, the imputation of specific values for these distributions, and the use of the 75% rule to make decisions about the presence or absence of heterogeneity. These practices may be consequences of psychometric meta-analysis origins in test validity research, and may resolve themselves as relevant data accumulate from other research domains.

### SUMMARY POINTS

- Psychometric meta-analysis attempts to correct for bias in study findings by adjusting for errors of measurement, restriction of range and other artifactual influences on effect sizes.
- Making these corrections typically produces effects that are higher in magnitude, and less variable across the set of studies in the meta-analysis, than would be the case if no corrections had been made.
- Psychometric meta-analysis emphasizes the distribution of true effects, rather than the overall mean effect. The inclusion of an estimate of the distribution of true effects should be adopted even by those who use conventional meta-analytic techniques.

### Further Reading

- Aguinis, H. (2001). Estimation of sampling variance of correlation in meta-analysis. *Personnel Psychology* 54: 569–590.
- Aguinis, H. & Pierce, C.A. (1998). Testing moderator variable hypotheses meta-analytically. *Journal of Management* 24: 577–592.
- Aguinis, H., Sturman, M., & Pierce, C.A. (2008). Comparison of three meta-analytic procedures for estimating moderating effects of categorical variables. *Organizational Research Methods* 11: 9–34.
- Bobko, P. & Roth, P.L. (2008). Psychometric accuracy and (the continuing need for) quality thinking in meta-analysis. *Organizational Research Methods* 11: 114–126.
- Crocker, L. & Algina, J. (1986). *Introduction to Classical and Modern Test Theory*. New York: Holt, Rinehart, & Winston.
- Hunter, J.E. & Schmidt, F.L. (2004). *Methods of Meta-analysis: Correcting Error and Bias in Research Findings* (2nd edn). Thousand Oaks, CA: Sage.
- Lord, F.M. & Novick, M.R. (1968). *Statistical Theories of Mental Test Scores*. Reading, MA: Addison-Wesley.
- Murphy, K.R. (2003). *Validity Generalization: A Critical Review*. Mahwah, NJ: Lawrence Erlbaum.
- Schmidt, F.L. & Hunter, J.E. (1977). Development of a general solution to the problem of validity generalization. *Journal of Applied Psychology* 62: 529–540.



# Meta-Analysis in Context



## Overview

Several chapters in this section address basic issues in meta-analysis. Chapter 45 addresses the basic question of when it makes sense to perform a meta-analysis, and Chapter 46 offers suggestions for how to report the results of a meta-analysis. The two share the theme that there is no single best way to perform (or report) a meta-analysis. Rather, there needs to be a match between the goals of a specific synthesis, the kinds of studies included, the methods applied, and the conclusions reported.

Chapter 47 describes a procedure called cumulative meta-analysis, which is a sequence of analyses performed with one study, then two, and so on, until all studies have been entered. This can be used to see how the body of evidence has shifted over time, or as a function of any moderator.

Chapter 48 is dedicated to criticisms of meta-analysis. There, we present a series of questions often raised by critics. We argue that some of the criticisms represent either a misunderstanding of the method, or poor applications of the method. Other criticisms represent problems that cannot easily be resolved, but we try to place these in context, by showing that these problems exist also in other kinds of reviews.



# When Does it Make Sense to Perform a Meta-Analysis?

---

### Introduction

Are the studies similar enough to combine?

Can I combine studies with different designs?

How many studies are enough to carry out a meta-analysis?

---

### INTRODUCTION

In the early days of meta-analysis (at least in its current incarnation) Robert Rosenthal was asked if it makes sense to perform a meta-analysis, given that the studies differ in various ways, and the analysis amounts to *combining apples and oranges*. Rosenthal answered that combining apples and oranges makes sense if your goal is to produce a fruit salad.

The goal of a meta-analysis is only rarely to synthesize data from a set of identical studies. Almost invariably, the goal is to broaden the base of studies in some way, expand the question, and study the pattern of answers. The question of whether it makes sense to perform a meta-analysis, and the question of what kinds of studies to include, must be asked and answered in the context of specific goals.

The ability to combine data from different studies to estimate the common effect (or mean effect) continues to be an important function of meta-analysis. However, it is not the only function. The goal of some syntheses will be to report the summary effect, but the goal of other syntheses will be to assess the dispersion as well as the mean effect, and the goal of others will be to focus on the dispersion exclusively.

For example, suppose that we are looking at the impact of a teaching intervention on student performance. Does it make sense to include studies that measured verbal skills and also studies that measured math skills? If our goal is to assess the impact on performance in general, then the answer is *Yes*. If our goal is to assess the impact on verbal skills alone, then the answer is *No*. Does it make sense to include studies that

enrolled middle-school students and also studies that enrolled high-school students? Again, the answer depends on the question being asked.

In some cases, however, the decisions are less clear cut. For example, does it make sense to include both randomized trials and observational studies in the same analysis? What about quasi-experimental studies? Is it acceptable to include studies that used independent groups and also studies that used matched designs? The answers to these and other questions will need to be decided in the context of the research question being addressed. Our goal here is to outline the kinds of issues that may arise and provide a context for making these kinds of decisions.

### ARE THE STUDIES SIMILAR ENOUGH TO COMBINE?

From a statistical perspective, there is no restriction on the similarity of studies based on the types of participants, interventions, or exposures. However, for the analysis to be meaningful, we need to pay careful consideration to the diversity of studies in these respects. For example, the research question might be *Does Drug A reduce the risk of heart attack as compared with a placebo; when used in a population of males, age 40–60 years, with no prior history of heart attack, and a cholesterol level of 250–300 on initial screening; where the dose was between 10–15 mg per day; in studies that followed patients for at least a year, with a drop-out rate no higher than five percent, where the randomization and blinding met specific criteria; and where patients were under the care of a primary care physicians for the duration of the study?* In this case, the criteria for including studies in the analysis would be very narrow, and the goal of the analysis would be to yield a more precise estimate (and more powerful test) of the effect than would be possible with any single study. This kind of meta-analysis might be planned by a pharmaceutical company as part of the approval process, and this approach is entirely legitimate.

In most meta-analyses, however, the inclusion criteria will be broader than this. It is an important feature of a meta-analysis that it may (and usually must) address a broader question than those addressed by the primary studies it includes. Thus a certain amount of diversity among the studies is not only inevitable but also desirable. A good meta-analysis will anticipate this diversity and will interpret the findings with attention to the dispersion of results across studies. To modify the prior example by relaxing some of the criteria, a *pragmatic* review of the effects of Drug A versus a placebo on the risk of heart attack might include both sexes, adults of any age with no prior history of heart attack, any cholesterol level; any dose of drug; in studies that followed patients for at least a year, with a drop-out rate no higher than twenty percent, where the randomization and blinding met specific criteria. The diversity of studies meeting these broader criteria may lead to heterogeneous results, and this heterogeneity needs to be recognized in the analysis and interpretation.

One approach to diversity is to apply the random-effects model and then address the diversity by reporting the expected range of true effects over the populations and interventions sampled. This could take the form of a prediction interval as explained

in Chapter 17. This is appropriate if the effects fall over a small range, so that the substantive implications of the finding are the same across the range.

With sufficient data, we can also explore the diversity across studies. For example, we could investigate how the effect of a drug (as compared with a placebo) depends on the sex of a patient. Assume that these studies reported outcomes for males and females separately. We now have the ability to compute a summary effect in each of these subgroups and to determine whether (and how) the effect is related to sex. If the effect is similar in both groups, then we can report that the effect is robust, something that was not possible with the more narrow criteria. If the effect varies (say, the drug is effective for males but not for females) then the meta-analysis may have yielded important information that, again, was not possible when all studies adhered to the more narrow criteria. Note, however, that for many meta-analyses there is insufficient power to do this reliably. There may also be problems of confounding.

This basic approach, that we can define the inclusion criteria narrowly and focus on the summary effect, or define the inclusion criteria more broadly and explore the dispersion, holds true for any meta-analysis. This idea can play itself out in various ways, and we explore some of them here.

## CAN I COMBINE STUDIES WITH DIFFERENT DESIGNS?

The appropriate types of study to include in a meta-analysis depend primarily on the type of question being addressed. For example, meta-analyses to evaluate the effect of an intervention will tend to seek randomized trials, in which interventions are assigned in an experimental fashion so that there are no important differences between those receiving and not receiving the intervention of interest. Meta-analyses to investigate the cause of a rare disease will tend to seek case-control studies, in which the past exposures of a collection of people with the disease are compared with those of a collection of people without the disease. Meta-analyses to examine the prevalence of a condition or a belief will tend to seek cross-sectional studies or surveys, in which a single group is examined and no within-study comparisons are made. And so on. Nevertheless, for any particular question there are typically several types of study that could yield a meaningful answer. A frequent question is whether studies with different designs can be combined in a metaanalysis.

### Randomized trials versus observational studies

Some have argued that systematic reviews on the effects of interventions should be limited to randomized controlled trials, since these are protected from internal bias by design, and should exclude observational studies, since the effect sizes in these are almost invariably affected by confounders (and the confounders may vary from one study to the next). In our opinion, this distinction is somewhat arbitrary. It suggests that we would be better off with a set of poor-quality randomized trials than with a set of high-quality observational studies (and leaves open the question of quasi-experimental

studies). The key distinction should not be the design of the studies but the extent to which the studies are able to yield an unbiased estimate of the effect size in question.

For example, suppose we wish to evaluate the effects of going to a support group to give up smoking. We might locate five studies in which smokers were recruited and then randomly assigned to either of two conditions (invitation to a support group, or a control intervention). Because the trials use random assignment, differences between groups are attributed to differences in the effects of the interventions. If we include these trials in a meta-analysis, we are able to obtain a more precise estimate of the effect than we could from any single trial, and this effect can be attributed to the treatment. However, since trials cannot *impose* an intervention on people, the effect is of being *invited* to the support group rather than, necessarily, of attending the support group. Furthermore, the types of smokers who volunteer to be randomized into a trial may not be the types of smokers who might volunteer to join a support group.

Alternatively, suppose that we locate five studies that compared the outcomes of smokers who had voluntarily joined a support group with others who had not. Because these studies are observational, any differences allow us to draw conclusions about what proportions of people are likely to be smoking after joining a support group or not joining a support group, but do not allow us to attribute these differences to the treatment itself. For instance, those who enrolled for treatment are likely to have been more motivated to stop smoking. If we include these observational studies in a meta-analysis we are able to obtain a more precise estimate of the difference than we could from any single study, but the interpretation of this difference is subject to the same limitations as that of the primary studies.

Does it make sense to include both these randomized trials and these observational studies in the same meta-analysis? The two kinds of studies are asking different questions. The randomized trial asks if there is a relationship between treatment and outcome when we control for all other factors, while the observational study asks if there is a relationship when we do not control for these factors. Furthermore, the *treatment* is different, in that the randomized trial evaluates the effect of the invitation, and the observational study collects information based on actual participation in the support group.

It would probably not make sense to compute a summary value across both kinds of studies. The meta-analyst should first decide which question is of greater interest. Unfortunately neither would seem to address the fundamental question of whether participating in the support group increases the likelihood of stopping smoking. As is often the case, the researcher must decide between asking the sub-optimal question (about invitations) with minimal bias (through randomization) or the right question (about participation) with likely bias (using observational studies). Most would argue that randomized trials do ask highly relevant questions, allowing important conclusions to be drawn about causality even if they do not fully reflect the way intervention would be applied on a day-to-day basis. Thus the majority of meta-analyses of interventions are restricted to randomized trials, at least in health care, where randomized trials have long been the established method of evaluation. Of course, some important effects of interventions, such as long-term or rare outcomes (especially harms) often

cannot be studied in randomized trials, so may need to be addressed using observational studies. We would generally recommend that randomized trials and observational studies be analyzed separately, though they might be put together if they do not disagree with each other and are believed to address a common question.

### **Studies that used independent groups, paired groups, clustered groups**

Suppose that some of the studies compared means for treatment versus control using two independent groups, others compared means using paired groups and others used cluster-randomized trials. There is no technical problem with combining data from the three kinds of studies, but we need to assume that the studies are functionally similar in all other important respects. On the one hand, studies that used different designs may differ from each other in substantive ways as well. On the other hand, these differences may be no more important than the difference between (say) studies that enrolled subjects in cities and others that enrolled subjects in rural areas. If we are looking at the impact of a vaccination, then the biological function is probably the same in all three kinds of studies. If we are looking at the impact of an educational intervention, then we would probably want to test this assumption rather than take it on faith.

### **Can I combine studies that report results in different ways?**

Meta-analysts frequently have to deal with results reported in different ways. Suppose we are looking at ways to increase the yield of grain, and are interested in whether a high dose of fertilizer works better than the standard dose. We might find studies that measure the impact of dose by randomizing different plots to receive one of the two doses, but which measure the outcome in different ways. Some studies might measure the average growth rate for the plants while others measure the yield after a certain number of weeks (and the timings might vary across studies). Some studies might measure the proportion of plants achieving a specific growth rate while others measure the time from application to production of a certain volume of grain. We might find further studies that apply a range of doses and examine the correlation between the dose and, for example, yield.

Even within studies investigating the same outcome, results can be reported in different ways. There are two types of variation here. First, different approaches to analysis could be used. For example, two studies might focus on the proportion of plants that fail under each dose of fertilizer, but one reports this as a ratio while another reports this as a difference in proportions. Second, even the same analysis can be reported using different statistics. For example, if several studies compare the mean yields between the two doses, some may report means with standard deviations, others means with a *p*-value, others differences in means with confidence intervals, and others *F* statistics from analysis of variance.

To what extent can all of these variations be combined in a meta-analysis? We address here only the statistical considerations, and assume that there is sound

rationale for combining the different outcome measures in the analysis. Note that we have described binary outcomes (proportion of failing plants), continuous outcomes using different measurement scales (growth rate, yield), survival outcomes (time to fruit) and correlational data (dose–yield). The list of possibilities is longer, and we do not attempt a comprehensive summary of all options.

When studies are addressing the same outcome, measured in the same way, using the same approach to analysis, but presenting results in different ways, then the only obstacles to meta-analysis are practical. If sufficient information is available to estimate the effect size of interest, then a meta-analysis is possible. For instance, means with standard deviations, means with a *p*-value, and differences in means with a confidence interval can all be used to estimate the difference in mean yield (providing, in the first two situations, that the sample sizes are known). These three also allow calculation of a standardized difference in means, as does a suitable *F* statistic in combination with sample size. Detailed discussions of such conversions are provided in Borenstein *et al.* (in preparation).

When studies are addressing the same outcome, measured in the same way, but using different approaches to analysis, then the possibility of a meta-analysis depends on both statistical and practical considerations. One important point is that all studies in a meta-analysis must use essentially the same index of treatment effect. For example, we cannot combine a risk difference with a risk ratio. Rather, we would need to use the summary data to compute the same index for all studies.

There are some indices that are similar, if not exactly the same, and judgments are required as to whether it is acceptable to combine them. One example is odds ratios and risk ratios. When the event is rare, then these are approximately equal and can readily be combined. As the event gets more common the two diverge and should not be combined. Other indices that are similar to risk ratios are hazard ratios and rate ratios. Some people decide these are similar enough to combine; others do not. The judgment of the meta-analyst in the context of the aims of the meta-analysis will be required to make such decisions on a case by case basis.

When studies are addressing the same outcome measured in different ways, or different outcomes altogether, then the suitability of a meta-analysis depends mainly on substantive considerations. The researcher will have to decide whether a combined analysis would have a meaningful interpretation. If so, then the above statistical and practical considerations apply. A further consideration is how different scales used for different outcomes are to be dealt with. The standard approach for continuous outcome measures is to analyze each study as a standardized mean difference, so that all studies share a common metric.

There is a useful class of indices that are, perhaps surprisingly, combinable under some simple transformations. In particular, formulas are available to convert standardized mean differences, odds ratios and correlations to a common metric (see Chapter 7). These kinds of conversions require some assumptions about the underlying nature of the data, and violations of these assumptions can have an impact on the validity of the process. Also, we must remember that studies which used dichotomous data may be different in some substantive ways than studies which used continuous data, and studies measuring correlations may be different from those

that compared two groups. As before, these are questions of degree rather than of qualitative differences among the studies.

## HOW MANY STUDIES ARE ENOUGH TO CARRY OUT A META-ANALYSIS?

If we are working with a fixed-effect model, then it makes sense to perform a meta-analysis as soon as we have two studies, since a summary based on two or more studies yields a more precise estimate of the true effect than either study alone. Importantly, we are not concerned with dispersion in the observed effects because this is assumed to reflect nothing more than sampling error. There might be a concern that by reporting a summary effect we are implying a level of certainty that is not warranted. In fact, though, the summary effect is qualified by a confidence interval that describes the uncertainty of the estimate. Additionally, research shows that if we fail to provide this information researchers will impose their own synthesis on the data, which will invariably be less accurate and more idiosyncratic than the value than we compute using known formulas.

In most cases, however, we should be working with the random-effects model, where the dispersion in effects is assumed to be real (at least in part). Unlike the fixed-effect analysis, where the estimate of the error is based on sampling theory (and therefore reliable), in a random-effects analysis, our estimate of the error may itself be unreliable. Specifically, when based on a small number of studies, the estimate of the between-studies variance ( $T^2$ ), may be substantially in error. The standard error of the summary effect is based (in part) on this value, and therefore, if we present a summary effect with confidence interval, not only is the point estimate likely to be wrong but the confidence interval may provide a false sense of assurance.

A separate problem is that in a random-effects analysis, our understanding of the dispersion affects not only our estimate of the summary effect but also the thrust of the analysis. In other words, if the effect is consistent across studies we would report that the effect is robust. By contrast, if the effect varies substantially from study to study we would want to consider the impact of the dispersion. The problem is that when we have only a few studies to work with, we may not know what the dispersion actually looks like.

This suggests that if the number of studies is small enough it might be better not to summarize them statistically. However many statisticians would argue that, when faced with a series of studies, people have an almost irresistible tendency to draw some summary conclusions from them. Experience has shown that seemingly intuitive *ad hoc* summaries (such as vote counting, Chapter 33) are also often highly misleading. This suggests that a statistical summary with known, but perhaps poor, properties (such as high uncertainty) may be superior to inviting an *ad hoc* summary with unknown properties.

In sum, when the number of studies is small, there are no really good options. As a starting point we would suggest reporting the usual statistics and then explaining the limitations as clearly as possible. This helps preclude the kinds of *ad hoc* analyses mentioned in the previous paragraph, and is an accurate representation of what we can do with limited data.

## SUMMARY POINTS

- The question of whether or not it makes sense to perform a meta-analysis is a question of matching the synthesis to the research question.
- If our goal is to report a summary effect, then the populations and interventions (and other variables) in the studies should match those in our target population. If our goal is to report on the dispersion of effects as a function of a covariate, then the synthesis must include the relevant studies and the analysis should focus on the differences in effects.
- Generally, we need to be aware of substantive differences among studies, but technical differences can be addressed in the analysis.
- A potentially serious problem exists when the synthesis is based on a small number of studies. Without sufficient numbers of studies we will have a problem estimating the between-studies variance, which has important implications for many aspects of the analysis.

## Further Reading

Ioannidis, J.P.A., Patsopoulos, N.A., Rothstein, H.R. (2008). Reasons or excuses for avoiding meta-analysis in forest plots. *BMJ* 336: 1413–1415.

# Reporting the Results of a Meta-Analysis

---

## Introduction

The computational model

Forest plots

Sensitivity analysis

---

## INTRODUCTION

Most of the issues that one would address when reporting the results of a meta-analysis are similar to those for reporting the results of a primary study. There are some unique issues as well, and we address those here.

As we have throughout this volume, we deal here only with issues related to the meta-analysis, and not to the full systematic review. For a broader perspective, especially for reviews in medicine, see the *Cochrane Handbook for Systematic Reviews of Interventions* (Higgins et al., 2019) and consider using some published Cochrane reviews or Campbell reviews as models for the full report.

## Are the effects consistent?

A recurring theme in this volume has been that the goal of a meta-analysis is to synthesize the data, which is not necessarily the same as computing a single summary effect. Similarly, the issues discussed in the report should match those that are important to the synthesis.

- If the effect sizes are consistent, then the focus of the report is likely to be on the summary effect, and the fact that the effect size is robust, in that it does not vary across the range of studies included in the analysis.
- If the effect sizes vary modestly from study to study we may still report the summary effect size, but will want to pay attention also to the dispersion in effects. Usually,

when we talk about whether or not the effects vary we are referring to substantive variation, so it would be helpful to report the range across which the true effects vary. If there are sufficient studies in the analysis, we will be able to estimate this range with reasonable precision. If this is not possible, we should acknowledge this limitation. Similarly, it would be important to report  $I^2$ , the proportion of total variance attributed to variance in true effects. This helps to place the observed dispersion in context (see Chapter 19).

- If the effect sizes vary substantially, the report might focus on the variance, with the summary effect being of less (or even no) importance.

In all cases we need to be careful about the meaning of dispersion. We need to distinguish between the case where the effects are shown to be homogeneous (on the one hand) and the case where we simply fail to reject the hypothesis of homogeneity (on the other). Also, one measure of heterogeneity ( $I^2$ ) tells us what proportion of the observed dispersion reflects differences in the true effect while others ( $T^2$ ,  $T$ , and the prediction interval) reflect the amount of heterogeneity on an absolute scale. While these measures tend to move in tandem, it is important to recognize that they are addressing two completely different aspects of heterogeneity, and it is important to use each appropriately in the report (see Chapter 17 on prediction intervals).

## THE COMPUTATIONAL MODEL

A report should state the computational model used in the analysis and explain why this model was selected. A common mistake is to use the fixed-effect model on the basis that there is no evidence of heterogeneity. As explained in Chapter 13, the decision to use one model or the other should depend on the nature of the studies, and not on the significance of this test.

## FOREST PLOTS

A recurring theme in this volume is the importance of interpreting statistics in context, and the forest plot helps to provide that context. The plot, as suggested by its appellation, allows the researcher to see both the forest and the trees (in the UK this would be the wood and the trees). We have used forest plots throughout this volume to illustrate various conceptual issues, precisely because the forest plot is an excellent vehicle for illustrating these issues. It can, and should, serve the same purpose in a report of a meta-analysis.

In the forest plot each study as well as the summary effect is depicted as a point estimate bounded by its confidence interval. It shows if the overall effect is based on many studies or a few, on studies that are precise or imprecise; whether the effects for all studies tend to line up, or whether they vary substantially from one study to the next. The plot puts a face on the statistics, helping to ensure that they will be interpreted properly, and highlighting anomalies, such as outliers, that require attention. The forest

plot is a compelling piece of information and easy to understand, even for people who do not work with meta-analysis on a regular basis.

There are several variants of the forest plot that appear in the literature. Our goal here is to sensitize the reader to some of these variants, and point out the advantages and disadvantages of each.

Consider Figures 46.1 and 46.2, which represent the same set of studies. In both, the study is represented by a point which is bounded by the confidence interval for the effect size in that study. In Figure 46.1 the point is a vertical line, while in Figure 46.2 the point is a box, proportional (in area) to that study's weight in the analysis.

Both plots use confidence intervals to track the precision, with a narrower interval reflecting better precision. What distinguishes between the two versions is the mechanism used to reflect the study's weight in the analysis, and the second version offers two advantages.

First, the boxes provide an important visual cue. In Figure 46.1 the only cue to the study weight is the width of the confidence interval, which requires careful attention and also is *inversely* proportional to the study weight. By contrast, in Figure 46.2 the studies with more weight are assigned *more* ink in the plot. The eye is naturally drawn to these studies, and we quickly get a sense of the relative impact of the different studies. In this example it is immediately apparent that the studies by Donat and by Young are dominant factors in the summary effect.

Second, the confidence interval (more precisely, the inverse of the squared standard error) is directly related to the weight only under the fixed-effect model. Under the random-effects model, the study weight is based also on the between-studies variance and may bear little relationship to the confidence interval.

The label *forest plot* also appears in the literature as *Forrest plot*. A paper in *BMJ* by Lewis and Clarke (2001) explains that this (incorrect) usage stems from a comment

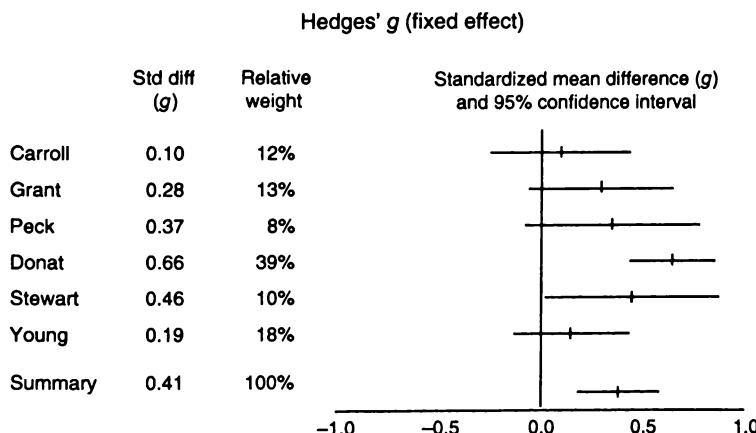


Figure 46.1 Forest plot using lines to represent the effect size.

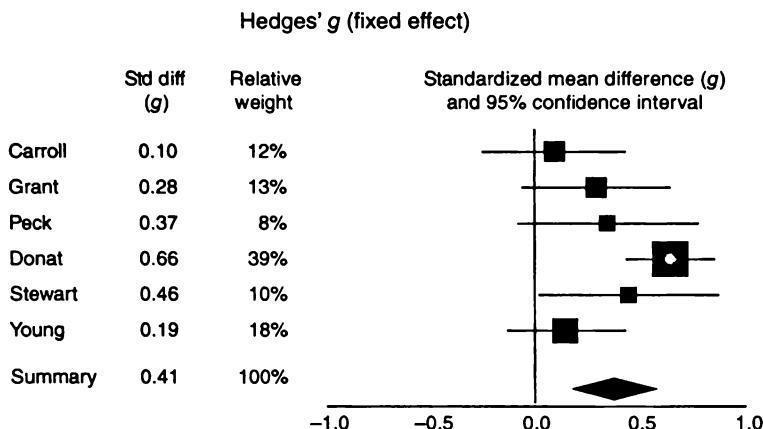


Figure 46.2 Forest plot using boxes to represent the effect size and relative weight.

made by Richard Peto who, as a joke, attributed the plot's invention to breast cancer researcher Pat Forrest.

## SENSITIVITY ANALYSIS

The issues addressed by a sensitivity analysis for a systematic review are similar to those that might be addressed by a sensitivity analysis for a primary study. That is, the focus is on the extent to which the results are (or are not) robust to assumptions and decisions that were made when carrying out the synthesis. The kinds of issues that need to be included in a sensitivity analysis will vary from one synthesis to the next. Our goal here is not to describe all of the possible sensitivity analyses that might be done in a meta-analysis, but to outline the kinds of issues that one might want to consider.

One kind of sensitivity analysis is concerned with the impact of decisions that lead to different data being used in the analysis. A common example of sensitivity analysis is to ask how results might have changed if different study inclusion rules had been used. This could be asked about studies classified on the basis of *a priori* criteria (for example, how would results have differed if we had included only randomized experiments instead of also including well-designed quasi-experiments). It could also be asked about studies identified as outliers (studies whose effects differ very substantially from the others). Here the question is whether the conclusions reached might differ substantially if a single study or a few studies were omitted.

Another kind of sensitivity analysis is concerned with the impact of the statistical methods used on the conclusions drawn from the analysis. For example one might ask whether the conclusions would have been different if a different effect size measure had been used (e.g. a risk ratio versus an odds ratio or an effect size using covariate adjusted means versus raw means). Alternatively, one might ask whether the

conclusions would be the same if fixed-effect versus random-effects methods had been used. We might also ask whether conclusions would be different if the analysis had adjusted for the effects of unreliability or restriction of range within individual studies.

Yet another kind of sensitivity analysis is concerned with how we addressed missing data. One situation is missing data on study characteristics that might be used formally (as in a moderator analysis) or informally (as a basis for grouping or describing studies). A very important form of missing data is the missing data on effect sizes that may result from incomplete reporting or selective reporting of statistical results within studies. When data are selectively reported in a way that is related to the magnitude of the effect size (e.g., when results are only reported when they are statistically significant), such missing data can have biasing effects similar to publication bias on entire studies. In either case, we need to ask how the results would have changed if we had dealt with missing data in another way.

Missing data are not limited just to study characteristics that are potential moderators nor to effect sizes. In some cases information needed to compute effect size estimates or their variances may not be reported and will be imputed in the meta-analysis (such as pretest–posttest correlations used to compute effect sizes). In this context, sensitivity analyses can be used to investigate whether the conclusions would differ substantially across a range of plausible imputed values.

In the next chapter we discuss cumulative analyses, which show how the summary effect and variance shift as studies are added to the analysis. This approach can also be used as part of a sensitivity analysis, for example by showing how our conclusions would (or would not) shift as new studies (perhaps representing a broader range of populations) are added.

### SUMMARY POINTS

- The questions being asked by the analysis, as well as the empirical findings, will help to shape the structure of the report, and this will vary from one analysis to the next. If the effect size is consistent across all studies in the analysis we are likely to focus on this effect and the fact that it is consistent. If the effect size varies somewhat we will want to estimate the amount of dispersion in true effects and consider the implications of this dispersion. If the effect size varies substantially, or if a goal of the analysis had been to explore expected variation in effect size, then the report would likely focus on the dispersion itself.
- The report of a meta-analysis should generally include a forest plot. This provides an intuitive sense of the data, and helps to ensure that the statistics will be interpreted in context.
- A sensitivity analysis is important to determine how robust the findings are. It would be important to know how the findings would shift if we changed the criteria for including studies, or if we changed some of the assumptions that we made when performing the analysis.

## Further Reading

- Lewis, S. & Clarke, M. (2001). Forest plots: trying to see the wood and the trees. *BMJ* 322: 1479–1480.
- Light, R.J. & Pillemer, D.B. (1984). *Summing Up: The Science of Reviewing Research*. Cambridge, MA: Harvard University Press.
- Light, R.J. & Pillemer, D.B. (1994). The visual presentation and interpretation of meta-analyses. In Cooper, H.M. & Hedges, L.V. (eds), *The Handbook of Research Synthesis*. New York: Russell Sage Foundation.

# Cumulative Meta-Analysis

---

### Introduction

Why perform a cumulative meta-analysis?

---

### INTRODUCTION

A cumulative meta-analysis is a meta-analysis that is performed first with one study, then with two studies, and so on, until all relevant studies have been included in the analysis. As such, a cumulative analysis *is not a different analytic method* than a standard analysis, but simply *a mechanism for displaying a series of separate analyses* in one table or plot. When the series are sorted into a sequence based on some factor, the display shows how our estimate of the effect size (and its precision) shifts as a function of this factor. When the studies are sorted chronologically, the display shows how the evidence accumulated, and how the conclusions may have shifted, over a period of time.

For example, consider the systematic review published by Lau *et al.* (1992) that looked at the impact of streptokinase in preventing death following a myocardial infarction. Streptokinase is a drug that has the potential to dissolve the blood clot that is causing a heart attack, and thus reduce the damage to heart muscle.

The systematic review synthesizes data from 33 studies that had been published over a period of 29 years. All the studies followed the same basic format, with patients who had suffered a myocardial infarction being assigned to either streptokinase or a placebo, and physicians recording the mortality rates in each group.

The standard meta-analysis is shown in Figure 47.1. Fletcher appears on the first row, with a risk ratio of 0.229, 95% confidence interval of 0.030 to 1.750. The *p*-value is 0.155, the sample size is 23, and the year is 1959. Dewar appears on the next row, with a risk ratio of 0.571, 95% confidence interval from 0.196 to 1.665. The *p*-value is 0.305, the sample size is 42, and the year is 1963. And so on for the remaining 31 studies.

The studies varied substantially in size, with five having fewer than 40 patients while one (GISSI-1 in 1986) enrolled 11,712 patients and one (ISIS-2 in 1988) enrolled

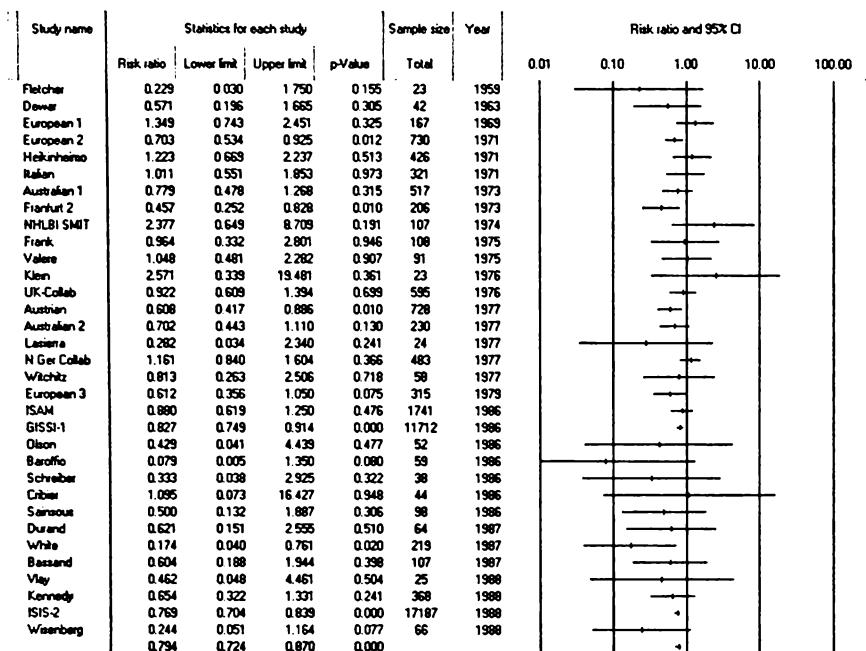


Figure 47.1 Impact of streptokinase on mortality – forest plot.

17,187 patients. Across all studies a total of 18,532 patients were assigned to treatment and 18,442 to control. The number of deaths in the two groups was 1892 versus 2375. The summary effect (using the random-effects model) is shown as a risk ratio of 0.794 with a 95% confidence interval from 0.724 to 0.870 and a *p*-value of 0.0000008.

The cumulative meta-analysis is shown in Figure 47.2. Here, we have the same 33 studies but the values on each row are not the statistics for that study. Rather, they are the summary values for a meta-analysis based on all studies up to and including that row. The line marked *Fletcher* is based only on Fletcher, and so is identical to the first line on the previous figure. The line marked *Dewar* shows the results of a meta-analysis based on Fletcher and Dewar. And so on. (Note that the scale of the forest plot has been changed.)

As one would expect, as we move down the plot the effect size tends to stabilize (as the volume of data accumulates, any new study is less likely to produce a sudden shift) and the confidence intervals tend to narrow (since the amount of data increases). The last study on the plot is Wisenberg. Since the analysis on this row includes data from all 33 studies, the statistics on this row are identical to those shown on the summary line. This also matches the summary line in Figure 47.1.

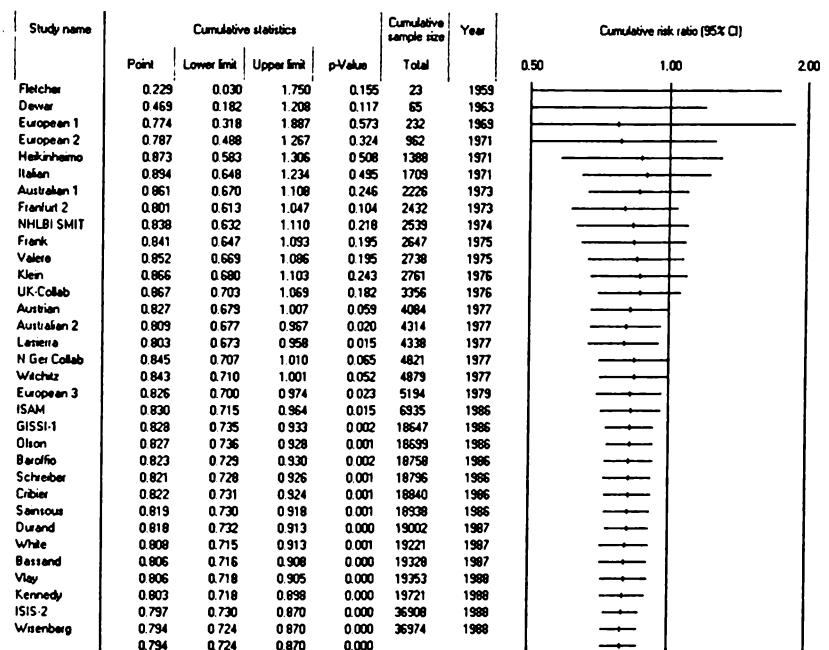


Figure 47.2 Impact of streptokinase on mortality – cumulative forest plot.

## WHY PERFORM A CUMULATIVE META-ANALYSIS?

### Cumulative meta-analysis as an educational tool

Lau *et al.* used the streptokinase analysis to show the potential impact of meta-analysis as part of the research process. They argued that if meta-analysis had been available to researchers several decades earlier, then the benefits of streptokinase could have been established as early as 1977. Had researchers performed a meta-analysis in 1977 using the studies prior to and including the Australian-2 study, they would have found that the risk ratio was 0.81, with a *p*-value of 0.020. Since meta-analysis was not yet recognized as a useful tool, researchers continued to perform additional studies (Lau *et al.*, 1992; Lau & Chalmers, 1995; Lau, Schmid, & Chalmers, 1995).

The studies published subsequent to the Australian-2 study enrolled a total of 32,660 patients, with approximately 50% assigned to placebo. In these studies there were 414 more deaths among the placebo patients than among the treated patients. Lau *et al.* argued that if meta-analysis had been conducted in 1977, then the efficacy of the treatment could have been established at that point and the subsequent trials could have been avoided. Not only would some of the patients who died on a placebo in these trials have been saved, but the drug would have become the standard of care and countless premature deaths worldwide could have been avoided.

One can argue with the specific numbers. In particular, when we repeatedly look at the cumulative data we may need to use a more conservative criterion for significance than 0.05 before deciding that the treatment is effective (see *Using a cumulative analysis prospectively* on page 411). Additionally, some trials that were published subsequent to the Australian-2 study were already underway when the Australian-2 study was published. Nevertheless, the basic argument is compelling, and captured the attention of the medical community. This cumulative analysis played an important role in gaining acceptance for meta-analysis as a useful mechanism for decision making.

The above should not be taken as a criticism of the people who performed the later studies. Meta-analysis was not widely accepted in the 1970s and 1980s and so was simply not an option at the time.

### To identify patterns in the data

While cumulative analyses are most often used to display the pattern of the evidence over time, the same technique can be used for other purposes as well. Rather than sort the data chronologically, we can sort it by any variable, and then display the pattern of effect sizes.

For example, assume that we have 100 studies that looked at the impact of homeopathic medicines, and we think that the effect is related to the quality of the blinding process. We anticipate that studies with complete blinding will show no effect, those with lower quality blinding will show a minor effect, those that blind only some people will show a larger effect, and so on. We could sort the studies based on the quality of the blinding (from high to low), and then perform a cumulative analysis. If our expectations were correct, the cumulative effect would initially be near zero, would increase as we moved to the next (lower) level of quality, and would increase some more at the next level.

Similarly, we could use cumulative analyses to display the possible impact of publication bias. The details will not be repeated here (this is covered in Chapter 35), but the problem being addressed is that the large studies are assumed to be unbiased, but the smaller studies may tend to over-estimate the effect size. We could perform a cumulative analysis, entering the larger studies at the top and adding the smaller studies at the bottom. If the effect was initially small when the large (nonbiased) studies were included, and then increased as the smaller studies were added, we would indeed be concerned that the effect size was related to sample size. A benefit of the cumulative analysis is that it displays not only *if* there is a shift in effect size, but also *the magnitude* of the shift.

### Display, not analysis

It is important to recognize that cumulative meta-analysis is a mechanism for display, rather than analysis. If we sort studies chronologically, we can see how the weight of the evidence has shifted over time. If we sort studies by the effectiveness of blinding,

we can see how the effect size shifts with the addition of poor-quality studies. If we sort studies by sample size, we can display the potential impact of publication bias.

These kinds of displays are compelling and can serve an important function. However, if our goal is actually to examine the relationship between a factor and effect size, then the appropriate analysis is a meta-regression, which looks at the relationship between each the study's effect size and the study's covariates (whether year of publication, or sample size, or something else).

### Using a cumulative analysis prospectively

As noted, the primary function of a cumulative analysis is to provide a mechanism for display. For example, when we sort the studies by year of publication and perform cumulative meta-analysis retrospectively, the analysis serves to provide us with historical context.

A very different situation emerges when the concept is applied prospectively (a process sometimes called *prospective* cumulative meta-analysis). In this case, a researcher performs a meta-analysis at time X using all the available data. If the cumulative effect is not definitive, then the analysis is repeated at time X + 1 with the addition of the next study (when that study becomes available). The process is repeated until such time as the results become definitive, at which time the process is stopped.

There is a serious problem with this use of cumulative analysis if the criterion for stopping is based on the analysis reaching a level of statistical significance. The problem is that the 0.05 (or any other) criterion only works as advertised when the data are subjected to a single statistical test. If the test is applied repeatedly, then (assuming that the null hypothesis is true) the likelihood of a false positive is 0.05 *for any given test*, but exceeds 0.05 when accumulated over all tests. This is analogous to a problem that arises in longitudinal studies where researchers follow a cohort of patients over time, look at the data periodically, and will stop the study if the *p*-value at any time crosses a given threshold.

There is an entire body of research on the best ways to allocate this risk and much of this can be applied to cumulative meta-analysis as well (Devereaux *et al.*, 2005; Pogue & Yusuf, 1998; Whitehead, 1997). For example, one solution used in longitudinal trials is to work with a more stringent criterion for stopping the trial, such as 0.01 rather than 0.05 (if 5 peeks are planned) so that the overall risk of a type I error is kept at an acceptable level.

At the same time, we need to recognize that researchers, clinicians and patients do not always have the luxury of waiting until enough studies have been completed before they make a decision. People who need to make a decision at a given point may need to base that decision on the evidence available at that point, and for this purpose a cumulative analysis of all available studies may be the best option. However, we can still continue to perform new studies and add to the cumulative evidence until the relevant questions have been fully addressed.

### SUMMARY POINTS

- Cumulative meta-analysis is a mechanism for displaying results from a series of separate analyses in one table or plot. The studies are typically sorted chronologically, which shows how the evidence has accumulated (and possibly how the results have shifted) over time. However, the studies may also be sorted by other variables, to show how the results shift as a function of some other factor (such as study quality).
- Cumulative analysis is a mechanism for display, not for analysis. If our goal is to test the hypothesis that the effect size has shifted over time, the correct approach would be to use subgroup analysis or meta-regression.
- A variant of cumulative analyses calls for study results to be added to a meta-analysis as each study is completed, with the analysis repeated every time the list of studies is updated. While this approach may provide the most up-to-date data for someone needing to make a decision about the utility of a treatment, if the plan is to stop adding studies when the analysis becomes definitive, then we need to adjust for the fact that we are having multiple looks at the data.

# Criticisms of Meta-Analysis

---

### Introduction

- One number cannot summarize a research field
  - The file drawer problem invalidates meta-analysis
  - Mixing apples and oranges
  - Garbage in, garbage out
  - Important studies are ignored
  - Meta-analysis can disagree with randomized trials
  - Meta-analyses are performed poorly
  - Is a narrative review better?
  - Concluding remarks
- 

### INTRODUCTION

While meta-analysis has been widely embraced by large segments of the research community, this point of view is not universal and people have voiced numerous criticisms of meta-analysis.

Some of these criticisms are worth mentioning for their creative use of metaphor. The first set of Cochrane reviews dealt with studies in neonatology, and one especially creative critic, cited by Mann (1990), called the reviewers *an obstetrical Baader Meinhof gang* (*obstetrical* being a reference to the field of research, and *Baader Meinhof gang* a reference to the terrorist group that operated in Europe during the 1970s and 1980s).

Others were more circumspect in their comments. Eysenck (1978) criticized a meta-analysis as *an exercise in mega-silliness*. Shapiro (1994) published a paper entitled *Meta-Analysis / Shmeta Analysis*. Feinstein (1995) wrote an editorial in which he referred to meta-analysis as ‘statistical alchemy for the 21st century’.

These critics share not only an affinity for allegory and alliteration but also a common set of concerns about meta-analysis. In this chapter we address the following criticisms that have been leveled at meta-analysis, as follows.

- One number cannot summarize a research field.
- The file drawer problem invalidates meta-analysis.
- Mixing apples and oranges.
- Garbage in, garbage out.
- Important studies are ignored.
- Meta-analysis can disagree with randomized trials.
- Meta-analyses are performed poorly.

After considering each of these questions in turn, we ask whether a traditional narrative review fares any better than a systematic review on these criticisms. And, we summarize the legitimate criticisms of meta-analysis that need to be considered whenever meta-analysis is applied.

## ONE NUMBER CANNOT SUMMARIZE A RESEARCH FIELD

### Criticism

A common criticism of meta-analysis is that the analysis focuses on the summary effect, and ignores the fact that the treatment effect may vary from study to study. Bailar (1997), for example, writes, ‘Any attempt to reduce results to a single value, with confidence bounds, is likely to lead to conclusions that are wrong, perhaps seriously so.’

### Response

In fact, the goal of a meta-analysis should be to *synthesize* the effect sizes, and not simply (or necessarily) to report a summary effect. If the effects are consistent, then the analysis shows that the effect is robust across the range of included studies. If there is modest dispersion, then this dispersion should serve to place the mean effect in context. If there is substantial dispersion, then the focus should shift from the summary effect to the dispersion itself. Researchers who report a summary effect and ignore heterogeneity are indeed missing the point of the synthesis.

## THE FILE DRAWER PROBLEM INVALIDATES META-ANALYSIS

### Criticism

While the meta-analysis will yield a mathematically sound synthesis of the studies included in the analysis, if these studies are a biased sample of all possible studies, then the mean effect reported by the meta-analysis will reflect this bias. Several lines of evidence show that studies finding relatively high treatment effects are more likely to be published than studies finding lower treatment effects. The latter, unpublished, research lies dormant in the researchers’ filing cabinets, and has led to the use of the term *file drawer problem* for meta-analysis.

## Response

Since published studies are more likely to be included in a meta-analysis than their unpublished counterparts, there is a legitimate concern that a meta-analysis may overestimate the true effect size.

Chapter 35 (entitled *Publication Bias*) explores this question in some detail. In that chapter we discuss methods to assess the likely amount of bias in any given meta-analysis, and to distinguish between analyses that can be considered robust to the impact of publication bias from those where the results should be considered suspect.

We must remember that publication bias is a problem for any kind of literature search. The problem exists for the clinician who searches a database to locate primary studies about the utility of a treatment. It exists for persons performing a narrative review. And, it exists for persons performing a meta-analysis. Publication bias has come to be identified with meta-analysis because meta-analysis has the goal of providing a more accurate synthesis than other methods, and so we are concerned with biases that will interfere with this goal. However, it would be a mistake to conclude that this bias is not a problem for the narrative review. There, it is simply easier to ignore.

## MIXING APPLES AND ORANGES

### Criticism

A common criticism of meta-analysis is that researchers combine different kinds of studies (*apples and oranges*) in the same analysis. The argument is that the summary effect will ignore possibly important differences across studies.

## Response

The studies that are brought together in a meta-analysis will inevitably differ in their characteristics, and the difficulty is deciding just how similar they need to be. The decision as to which studies should be included is always a judgment, and people will have different opinions on the appropriateness of combining results across studies. Some meta-analysts may make questionable judgments, and some critics may make unreasonable demands on similarity.

We need to remember that meta-analyses almost always, by their very nature, address broader questions than individual studies. Hence a meta-analysis may be thought of as asking a question about fruit, for which both apples and oranges (and indeed pears and melons) contribute valuable information. One of the strengths of meta-analysis is that the consistency, and hence generalizability, of findings from one type of study to the next can be assessed formally.

Of course, we always need to remember that we are dealing with different kinds of fruit, and to anticipate that effects may vary from one kind to the other. It is a further strength of meta-analysis that these differences, if identified, can be investigated formally. Assume, for example, that a treatment is very effective for patients with acute

symptoms but has no effect for patients with chronic symptoms. If we were to combine data from studies that used both types of patients, and conclude that the treatment was modestly effective (on average), this conclusion would not be accurate for either kind of patient. If we were to restrict our attention to studies in only patients with acute symptoms, or only patients with chronic symptoms, we could report how the treatment worked with one type of patient, but could only speculate about how it would have worked with the other type. By contrast, a meta-analysis that includes data for both types of patients may allow us to address this question empirically.

## GARBAGE IN, GARBAGE OUT

### Criticism

The often-heard metaphor *garbage in, garbage out* refers to the notion that if a meta-analysis includes many low-quality studies, then fundamental errors in the primary studies will be carried over to the meta-analysis, where the errors may be harder to identify.

### Response

Rather than thinking of meta-analysis as a process of *garbage in, garbage out* we can think of it as a process of waste management. A systematic review or meta-analysis will always have a set of inclusion criteria and these should include criteria based on the quality of the study. For trials, we may decide to limit the studies to those that use random assignment, or a placebo control. For observational studies we may decide to limit the studies to those where confounders were adequately addressed in the design or analysis. And so on. In fact, it is common in a systematic review to start with a large pool of studies and end with a much smaller set of studies after all inclusion/exclusion criteria are applied.

Nevertheless, the studies that do make it as far as a meta-analysis are unlikely to be perfect, and close attention should be paid to the possibility of bias due to study limitations. A meta-analysis of a collection of studies that is each biased in the same direction will suffer from the same bias and have higher precision. In this case, performing a meta-analysis can indeed be more dangerous than not performing one.

However, as noted in the response to the previous criticism about *apples and oranges*, a strength of meta-analysis is the ability to investigate whether variation in characteristics of studies is related to the size of the effect. Suppose that ten studies used an acceptable method to randomize patients while another ten used a questionable method. In the analysis we can compare the effect size in these two subgroups, and determine whether or not the effect size actually differs between the two. Note that such analyses (those comparing effects in different subgroups) can have very low power so need to be interpreted carefully, especially when there are not many studies within subgroups.

## IMPORTANT STUDIES ARE IGNORED

### Criticism

Whereas the *garbage in, garbage out* problem relates to the inclusion of studies that perhaps should not be included, a common complementary criticism is that important studies were left out. The criticism is often leveled by people who are uncomfortable with the findings of a meta-analysis. For example, a meta-analysis to assess the effects of antioxidant supplements (beta-carotene, vitamin A, vitamin C, vitamin E, and selenium) on overall mortality was met with accusations on the website of the Linus Pauling Institute (Oregon State University) that in this ‘flawed analysis of flawed data’ the authors looked at 815 human clinical trials of antioxidant supplements, but only 68 were included in the meta-analysis.

### Response

We have explained that systematic reviews and meta-analyses require explicit mechanisms for deciding which studies to include and which ones to exclude. These eligibility criteria are determined by a combination of considerations of relevance and considerations of bias, and are typically decided before the search for studies is implemented. Studies should be sufficiently similar to yield results that can be interpreted, and sufficiently free of bias to yield results that can be believed. For both purposes, judgments are required, and not all meta-analysts or readers would reach the same judgments on each occasion. Importantly, in meta-analysis the criteria are transparent and are described as part of the report.

## META-ANALYSIS CAN DISAGREE WITH RANDOMIZED TRIALS

### Criticism

LeLorier *et al.* (1997) published a paper in which they pointed out that meta-analyses sometimes yield different results than large-scale randomized trials. Specifically, they located cases in the medical literature where someone had performed a meta-analysis, and someone else subsequently performed a large-scale randomized trial that addressed the same question (e.g. *Does the treatment work?*). The authors reported that the results of the meta-analysis and the randomized trial *matched* (both were statistically significant, or neither was statistically significant) in about 66% of cases, but did not match (one was statistically significant but the other was not) in the remaining 34%. Since randomized trials are generally accepted as the gold standard they conclude that some 34% of these meta-analyses were wrong, and that meta-analyses in general cannot be trusted.

### Response

There are both technical and conceptual flaws in this criticism. The technical flaws relate to the question of what we mean by *matching*, and the authors’ decision to define

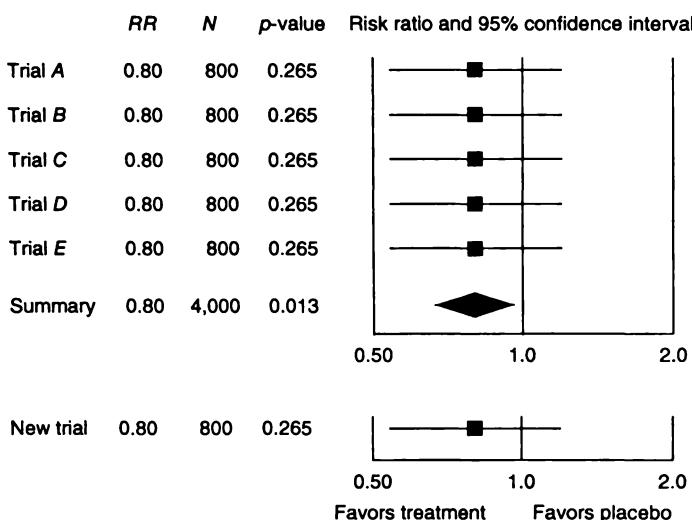


Figure 48.1 Forest plot of five fictional studies and a new trial (consistent effects).

matching as both studies being (or not being) statistically significant. The discussion that follows draws in part on comments by Ioannidis *et al.* (1998), Lelorier *et al.* (1997, 536–543) and others (see further readings at the end of this chapter).

Consider Figure 48.1, which shows a meta-analysis of five randomized controlled trials (RCTs) at the top, and a subsequent large-scale randomized trial at the bottom.

In this fictional example the five studies in the meta-analysis each showed precisely the same effect, an odds ratio of 0.80. The summary effect in the meta-analysis is (it follows) an odds ratio of 0.80. And, the subsequent study showed the same effect, an odds ratio of 0.80.

The only difference between the summary effect in the meta-analysis and the effect in the subsequent study is that the former is reported with greater precision (since it is based on more data) and therefore yields a *p*-value under 0.05. By the LeLorier criterion these two conclusions would be seen as conflicting, when in fact they have the identical effect size.

Additionally, LeLorier concludes that in the face of this conflict the single randomized trial is correct and the meta-analysis is wrong. In fact, though, it is the meta-analysis, which incorporates data from five randomized trials rather than one, that has the more powerful position. (What would happen if we performed a new meta-analysis which incorporated the most recent randomized trial? Would LeLorier now see this new meta-analysis as flawed?) In fact, the real issue is not that a meta-analysis disagrees with a randomized trial, but that randomized trials disagree with each other.

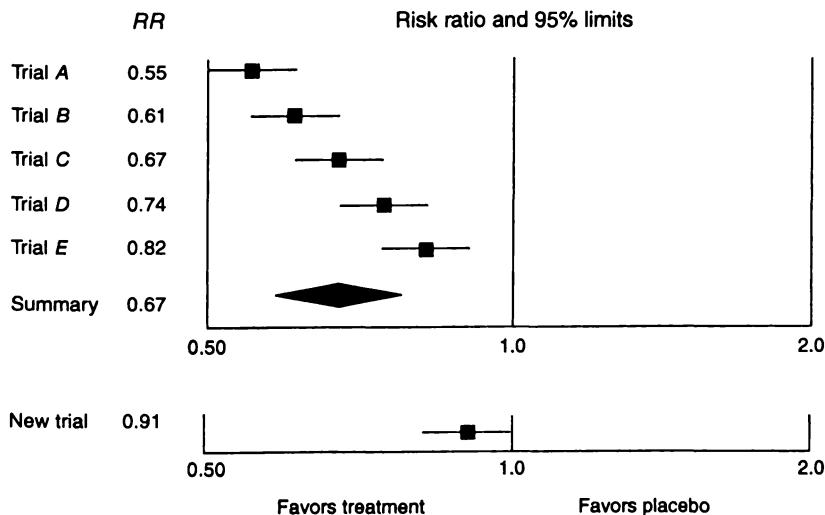
At a meeting of The Cochrane Collaboration in Baltimore (1996), a plenary speaker made the same argument being made by LeLorier *et al.* (that meta-analyses sometimes

yield different results than randomized trials) and, like the paper, cited the statistic that roughly a third of meta-analyses fail to match the *comparable* randomized trial. A distinguished member of the audience, Harris Cooper, asked the speaker if he knew what percentage of randomized trials fail to match the next randomized trial on the same topic. It turns out that the percentage is roughly a third.

However, to move on to a more interesting question, let's assume that the results from a meta-analysis and a randomized trial really do differ. Suppose that the meta-analysis yields a risk ratio of 0.67 (with a 95% confidence interval of 0.84 to 0.77) while the new trial yields a risk ratio of 0.91 (0.82 to 1.0). According to the meta-analysis the treatment reduces the risk by at least 23%, while the new trial says that its impact is no more than 18%.

In this case the effect is *different* in the two analyses, but that does not mean that one is wrong and the other is right. Rather, it behooves us to ask why the two results should differ, much as we would if we had two large-scale randomized trials with significantly different results. Often, it will turn out that the different analyses either were asking different questions or differed in some important way. A careful examination of the differences in method, patient population, and so on, may help to uncover the source of the difference.

Consider the following scenario, depicted in Figure 48.2. A new compound is introduced, which is meant to minimize neurological damage in stroke patients. In 1990, the compound is tested in a randomized trial involving patients with a very poor prognosis, and yields a risk ratio of 0.55. Based on these encouraging results, in 1994 it is tested in patients with a somewhat better prognosis. Since the patients in this group are



**Figure 48.2** Forest plot of five fictional studies and a new trial (heterogeneous effects).

more likely to recover without treatment, the impact of the drug is less pronounced, and the risk ratio is 0.61. By 1998 the drug is being tested with all patients, and the risk ratio is 0.82. These are the studies included in the meta-analysis. The new trial is performed using a relatively healthy population and (following the trend seen in the meta-analysis) yields a risk ratio of 0.91.

If one were to report a mean effect of 0.67 for the meta-analysis versus 0.91 for the new trial there would indeed be a problem. But, as we have emphasized throughout this volume, the meta-analysis should focus on the dispersion in effects and try to identify the reason for the dispersion. In this example, using either health status or study year as a covariate we can explain the pattern of the effects, and would have predicted that the effect size in the new study would fall where it did. For a comprehensive discussion of the replication issue, see Hedges and Schauer (2019).

## META-ANALYSES ARE PERFORMED POORLY

### Criticism

John C. Bailar, in an editorial for the *New England Journal of Medicine* (Bailar, 1997), writes that mistakes such as those outlined in the prior criticisms are common in meta-analysis. He argues that a meta-analysis is inherently so complicated that mistakes by the persons performing the analysis are all but inevitable. He also argues that journal editors are unlikely to uncover all of these mistakes.

### Response

The specific points made by Bailar about problems with meta-analysis are entirely reasonable. He is correct that many meta-analyses contain errors, some of them important ones. His list of potential (and common) problems can serve as a bullet list of mistakes to avoid when performing a meta-analysis.

However, the mistakes cited by Bailar are flaws in the application of the method, rather than problems with the method itself. Many primary studies suffer from flaws in the design, analyses, and conclusions. In fact, some serious kinds of problems are endemic in the literature. The response of the research community is to locate these flaws, consider their impact for the study in question, and (hopefully) take steps to avoid similar mistakes in the future. In the case of meta-analysis, as in the case of primary studies, we cannot condemn a method because some people have used that method improperly. As Bob Abelson once remarked in a related context, 'Think of all the things that people abuse. There are college educations. And oboes.'

## IS A NARRATIVE REVIEW BETTER?

In his editorial Bailar concludes that, until such time as the quality of meta-analyses is improved, he would prefer to work with the traditional narrative reviews: 'I still prefer conventional narrative reviews of the literature, a type of summary familiar to readers of the countless review articles on important medical issues.'

We disagree with the conclusion that narrative reviews are preferable to systematic reviews, and that meta-analyses should be avoided. The narrative review suffers from every one of the problems cited for the systematic review. The only difference is that, in the narrative review, these problems are less obvious. For example:

- The process of determining which studies to include in the systematic review or meta-analysis is difficult and prone to error. But at least there is a set of criteria for determining which studies to include. If the narrative review also has such criteria, then it is subject to the same kinds of error. If not, then we have no way of knowing how studies are being selected, which only compounds the problem.
- Meta-analyses can be affected by publication bias. But the same biases exist in the material upon which narrative reviews are based. Meta-analysis offers a means to investigate the likelihood of these biases and their potential impact on the results.
- Meta-analyses may be based on low-quality primary research. But a good systematic review includes a careful assessment of the included studies with regard to their quality or risk of bias, and meta-analytic methods enable formal examination of the potential impact of these biases. A narrative reviewer may discount a study because of a belief that the results are suspect for some reason. However, a limitation can be found for virtually any study, so in the absence of a systematic quality assessment of every study, a narrative reviewer is free to be suspect about any study's results and to lay the blame on one or more of its limitations.
- The weighting scheme in a meta-analysis may give a lot (or little) weight to specific studies in ways that may appear inappropriate. But in a meta-analysis the weights reflect specific goals (to minimize the variance, or to reflect the range of effects) and the weighting scheme is detailed as part of the report, so a reader is able to agree or disagree with it. By contrast, in the case of a narrative review, the reviewer assigns *weights* to studies based on criteria that he or she does not communicate, and may not even be able to fully articulate. Here, the problem involves not only the relative weights assigned to small or large studies. It extends also to the propensity of one reviewer to focus on effect sizes, and of another to focus on (and possibly be misled by) significance tests.
- Some meta-analyses focus on the summary effect and ignore the pattern of dispersion in the results. To ignore the dispersion is clearly a mistake both in a narrative review and in a meta-analysis. However, meta-analysis provides a full complement of tools to assess the pattern of dispersion, and possibly to explain it as a function of study-level covariates. By contrast, it would be an almost impossible task for a narrative reviewer to accurately assess the pattern of dispersion, or to understand its relationship to other variables.
- In support of the narrative review, Bailer cites the role of the expert with substantive knowledge of the field, who can identify flaws in specific studies, or the presence of potentially important moderator variables. However, this is not an advantage of the narrative review, since the expert is expected to play the same role in a meta-analysis. Steve Goodman (1991) wrote, 'The best meta-analyses knit clinical insight with quantitative results in a way that enhances both. They should combine the careful thought and synthesis of a good review with the scientific rigor of a good experiment.'

## CONCLUDING REMARKS

Most of the criticisms raised in this chapter point to problems with meta-analysis, and make the implicit argument that the problem would go away if we dispensed with the meta-analysis and performed a narrative review. We have argued that these problems exist also for the narrative review, and that the key advantage of the systematic approach of a meta-analysis is that all steps are clearly described so that the process is transparent.

Is meta-analysis so difficult that the method should be abandoned, as some have suggested? Our answer is obviously that it is not. Most of the criticisms raised deal with the application of the method, rather than with the method itself. What we should do is take the valid criticisms seriously and protect against them in planned analyses and by thoughtful interpretation of results.

Steven Goodman, in his editorial for *Annals of Internal Medicine* (1991) writes,

Regardless of the summary number, meta-analysis should shed light on why trial results differ; raise research and editorial standards by calling attention to the strengths and weaknesses of the body of research in an area; and give the practitioner an objective view of the research literature, unaffected by the sometimes distorting lens of individual experience and personal preference that can affect a less structured review.

### SUMMARY POINTS

- Meta-analyses are sometimes criticized for a number of flaws, and critics have argued that narrative reviews provide a better solution.
- Some of these flaws, such as the idea that we cannot summarize a body of data in a single number, are based on misunderstandings of meta-analysis.
- Many of the flaws (such as ignoring dispersion in effect sizes) reflect problems in the way that meta-analysis is used, rather than problems in the method itself.
- Other flaws (such as publication bias) are a problem for meta-analysis. However, the suggestion that these problems do not exist in narrative reviews is wrong. These problems exist for narrative reviews as well, but are simply easier to ignore since those reviews lack a clear structure.

### Further Reading

- Bailar, J.C. (1995). The practice of meta-analysis. *J Clin Epidemiol* 48: 149–157.
- Bailar, J.C. (1997). The promise and problems of meta-analysis. *New Engl J Med* 337: 559–561.
- Boden, W.E. (1992). Meta-analysis in clinical trials reporting: has a tool become a weapon? *Am J Cardiol* 69: 681–686.
- Egger, M., & Davey Smith, G. (1998). Bias in location and selection of studies. *BMJ* 316: 61–66.
- Eysenck, H.J. (1978). An exercise in mega-silliness. *Am Psychol* 33: 517.
- Hedges, L.V., & Schauer, J.M. (2019). More Than One Replication Study Is Needed for Unambiguous Tests of Replication. *Journal of Educational and Behavioral Statistics*, 44(5): 543–570.

- Lau, J., Ioannidis, J.P., Terrin, N., Schmid, C.H., & Olkin, I. (2006). The case of the misleading funnel plot. *BMJ* 333: 597–600.
- LeLorier, J., Gregoire, G., Benhaddad, A., Lapierre, J., & Derderian, F. (1997). Discrepancies between meta-analyses and subsequent large randomized, controlled trials. *N Engl J Med* 337: 536–543.

Responses to Lelorier *et al.*

- Bent, S., Kerlikowske, K., & Grady, D. (1998). *NEJM* 338(1): 60.
- Imperiale, T.F. (1998). *NEJM* 338(1): 61.
- Ioannidis, J.P., Cappelleri, J.C., & Lau, J. (1998). *NEJM* 338(1): 59.
- Khan, S., Williamson, P., & Sutton, R. (1998). *NEJM* 338(1): 60–61.
- LeLorier, J., & Gregoire, G. (1998). *NEJM* 338(1): 61–62.
- Song, F. J., & Sheldon, T. A. (1998). *NEJM* 338(1): 60.
- Stewart, L. A., Parmar, M. K., & Tierney, J. F. (1998). *NEJM* 338(1): 61

Sharpe, D. (1997). Of apples and oranges, file drawers and garbage: why validity issues in meta-analysis will not go away. *Clin Psychol Rev* 17: 881–901.

Thompson, S.G & Pocock, S. J. (1991). Can meta-analysis be trusted? *Lancet* 338: 1127–1130.



# Comprehensive Meta-Analysis Software

---

Introduction  
Features in CMA  
Teaching elements  
Documentation  
Availability  
Acknowledgments  
Motivating example  
Data entry  
Basic analysis  
What is the *average effect size*?  
How much does the effect size vary?  
Plot showing distribution of effects  
High-resolution plot  
Subgroup analysis  
Meta-regression  
Publication bias  
Explaining results

---

### INTRODUCTION

The screenshots in this volume are from the software Comprehensive Meta-Analysis (CMA). In this chapter, we provide an overview of this software and show how to use it to implement the ideas outlined in prior chapters. The same approach could be used with any other program as well. Our goal in this chapter is to provide a sense for the look-and-feel of the program. For the reader who would like to carry out the analyses, a step-by-step PDF is available on the book's website ([www.Introduction-to-Meta-Analysis.com](http://www.Introduction-to-Meta-Analysis.com)).

Full disclosure. The authors of this volume are also the developers of CMA, and some have a financial interest in the program. This chapter is a slightly modified

version of a chapter published in the text *Common Mistakes in Meta-Analysis – and How to Avoid Them* (Borenstein, 2019).

Comprehensive Meta-Analysis (CMA) is a computer program for meta-analysis that was developed with funding from the National Institutes of Health in the United States. The program was initially released in 2000 and has been updated on a regular basis since then.

CMA features a spreadsheet view and a menu-driven interface. As such, it allows a researcher to enter data and perform a simple analysis in a matter of minutes. At the same time, it offers a wide array of advanced features, including the ability to compare the effect size in subgroups of studies, to run meta-regression, to estimate the potential impact of publication bias, and to produce high-resolution plots. The program is designed to work with studies that compare an outcome in two groups or that estimate an outcome in one group. It is not intended for network meta-analyses nor for meta-analyses of diagnostic test accuracy.

## FEATURES IN CMA

The first step in conducting a meta-analysis is to compute an effect size and variance for each study. Many programs will perform this computation automatically when the data are in the form of  $2 \times 2$  tables or in the form of means and standard deviations for each group, but not for more complex data formats. By contrast, CMA will allow the user to enter data in more than one hundred formats. For example, the user can enter data as events and N in each group; or as an odds ratio and its confidence interval; or as a log risk ratio and its standard error. Or, the user can enter means and standard deviations for two independent groups; or the pre and post scores for a pre/post study; or the *p*-value from a *t*-test for two independent groups; and so on. Critically, the user may use a different format for each study. Thus, if one study reports means and standard deviations for two independent groups, a second reports a *p*-value based on two independent groups, and a third reports pre-scores and post-scores for a pre/post study, the user may enter data for each study in its own format. The program will apply the appropriate formula for each format to compute the effect size and its variance, and then include all the effects in the analysis.

The program allows the user to select from an array of effect size indices, including the odds ratio, risk ratio, risk difference, mean difference, standardized mean difference, correlation, hazard ratio, and prevalence, among others.

In the motivating example, each study provides one row of data. CMA also allows for the possibility that some (or all) studies will provide more than one row of data. There is an option for studies to report data for two or more outcomes, based on the same subjects. In the analysis, we could elect to look at either outcome alone. Or, we could tell the program to create a synthetic outcome which incorporates both measures, taking into account the fact that the two outcomes are not independent of each other.

Similarly, we can enter data for an outcome recorded at two or more time-points, which allows us to assess the impact at each time-point, and to see whether the effect

size changes over time. We can enter data for two or more independent subgroups within studies and then run the analysis using either subgroup or study as the unit of analysis. Finally, we can enter data for studies that employed one control group and multiple treatment groups.

## TEACHING ELEMENTS

The program incorporates a number of features intended to make the computations as transparent as possible. On the data-entry screen, the user enters summary data and the program displays the effect size and its variance. Double-click on the computed values, and the program will show how those values were computed. On the analysis screen, there is also a tab labeled 'Calculations', which opens a window onto the calculations.

## DOCUMENTATION

A manual is installed with the program. Each module in the program features an interactive guide that will walk the user through that module. Additionally, the website offers an array of PDFs and videos that show how to enter data, run the analysis, and then interpret the output. In each case, we also discuss how to report the data. The program's algorithms are discussed in this volume.

## AVAILABILITY

The program's website is [www.Meta-Analysis.com](http://www.Meta-Analysis.com). The program may be downloaded and run for free as a trial, and the website lists rates for licenses. There are discounts available for nonprofit institutions and for students. The program is free for short-term workshops in meta-analysis and is available at a discount for semester-length classes in meta-analysis.

## ACKNOWLEDGMENTS

Development of the program was funded by the National Institutes of Health in the United States under the following grants: MH052969 (Computer program for meta-analysis in mental health), AG021360 (Combining data types in meta-analysis), AG020052 (Publication bias in meta-analysis for mental health), AG024771 (Software for meta-regression), DA019280 (Forest plots for meta-analysis), AG029029 (Software for meta-analysis of diagnostic tests), and DA029351 (Software for meta-analysis with correlated outcomes). As a matter of policy, NIH does not endorse any product or software.

The program was developed by Michael Borenstein, Larry Hedges, Julian Higgins, and Hannah Rothstein. We gratefully acknowledge the contributions of Doug Altman, Betsy Becker, Jesse Berlin, Michael Brannick, Harris Cooper, Kay Dickersin,

Sue Duval, Roger Harbord, John Ioannidis, Jeff Valentine, Spyros Konstantopoulos, Mark Lipsey, Mike McDaniel, Fred Oswald, Terri Pigott, David Rindskopf, Stephen Senn, Will Shadish, Jonathan Sterne, Alex Sutton, Steven Tarlow, Thomas Trikalinos, Jack Vevea, Vish Viswesvaran, and David Wilson.

## MOTIVATING EXAMPLE

To illustrate the program, we will use a meta-analysis of 17 studies that assessed the utility of methylphenidate for treating adults with attention deficit hyperactivity disorder (ADHD). In each study, patients who had been diagnosed with ADHD were randomly assigned to either methylphenidate or a placebo and then tested on a scale intended to assess cognitive function (Castells *et al.*, 2011).

The effect size index is the standardized mean difference ( $d$ ). In this context, a standardized mean difference of 0.20 would be considered trivial – this is a difference that shows up on the tests, but the patient might not be aware of any change. A standardized mean difference of 0.50 would be considered moderate – the patient would recognize that they were doing better than usual, and coworkers might be aware of a change. A standardized mean difference of 0.80 would be considered large – the patient would feel great, and the difference would be obvious enough that others might remark on it.

## DATA ENTRY

Figure 49.1 shows the data-entry screen. For each study, enter the study name into column [A] and the summary data into the columns labeled [B]. The program displays

Study name	Std diff in means	Standard error	Group-A N (Optional)	Group-B N (Optional)	Effect direction	Std diff in means	Std Err	Variance	Dose	SUD
1 Levin a	-0.260	0.280			Auto	-0.260	0.280	0.078	60.0 Y	
2 Levin b	0.060	0.200			Auto	0.060	0.200	0.040	50.0 N	
3 Tenerbaum	0.070	0.290			Auto	0.070	0.290	0.084	45.0 N	
4 Carpenter	0.300	0.330			Auto	0.300	0.330	0.109	45.0 Y	
5 Guabao	0.310	0.510			Auto	0.310	0.510	0.260	48.7 N	
6 Medon	0.420	0.120			Auto	0.420	0.120	0.014	42.0 N	
7 Rosler	0.450	0.130			Auto	0.450	0.130	0.017	41.2 N	
8 Spencer c	0.510	0.160			Auto	0.510	0.160	0.026	29.8 N	
9 Adler	0.530	0.140			Auto	0.530	0.140	0.020	67.7 N	
10 Jan	0.540	0.240			Auto	0.540	0.240	0.058	56.8 N	
11 Wende	0.570	0.250			Auto	0.570	0.250	0.063	43.2 N	
12 Beutler	0.630	0.290			Auto	0.630	0.290	0.084	45.0 N	
13 Schubert	0.700	0.300			Auto	0.700	0.300	0.090	78.8 Y	
14 Berdeman	0.720	0.190			Auto	0.720	0.190	0.036	90.9 N	
15 Reinherz	0.830	0.360			Auto	0.830	0.260	0.068	64.0 N	
16 Spencer a	1.010	0.310			Auto	1.010	0.310	0.096	66.5 N	
17 Spencer b	1.300	0.280			Auto	1.300	0.280	0.078	82.0 N	
18										

Figure 49.1 Data-entry screen in CMA.

the standardized mean difference, standard error, and variance in the columns labeled [C]. We have also entered data for a series of moderator variables in the columns labeled [D], including the dose of methylphenidate (Dose), and whether the study enrolled patients who were abusing drugs (SUD).

In this example, the summary data entered for each study were the standardized mean difference and its standard error, because these are the data that had been reported. However, the user may elect to enter data in more than 100 formats. Similarly, we have elected to display the standardized mean difference, but the program will compute and display a wide array of effect size indices. The data may be entered directly into CMA or copied from another program such as Excel™.

To run the analysis, click [Run Analyses] on the toolbar.

## BASIC ANALYSIS

Figure 49.2 shows the analysis screen. A tab at the bottom [E] may be used to switch between fixed-effect and random-effects meta-analyses. The fixed-effect model is appropriate when the intended inference is limited to the studies in the analysis. The random-effects model is appropriate when the intended inference is to the universe of comparable studies. In this example, we intend to generalize the results to the universe of comparable studies and have selected the random-effects model (see Part 3).

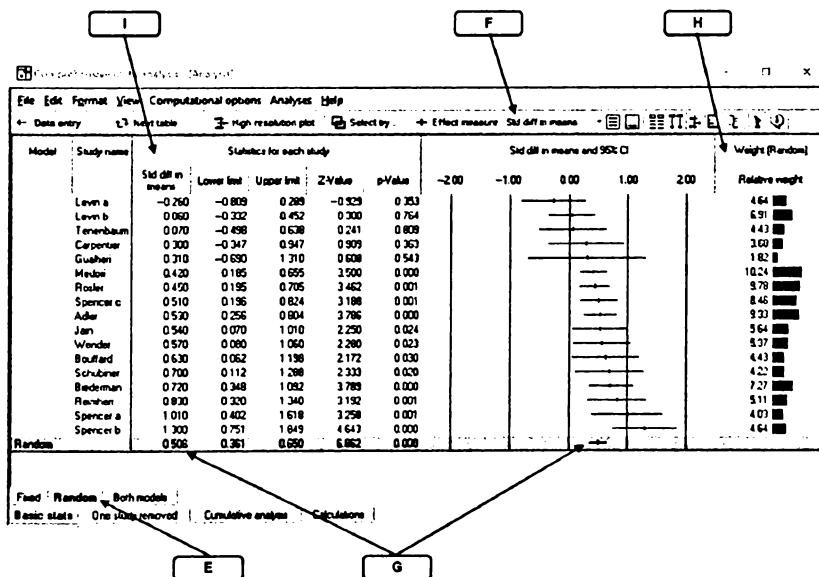


Figure 49.2 Basic analysis screen in CMA.

The toolbar [F] allows us to select the effect size index. Here, we have selected 'Standardized difference in means' (same as standardized mean difference).

## WHAT IS THE AVERAGE EFFECT SIZE?

The program [G] displays the average effect size as 0.506 and the confidence interval as 0.361 to 0.650. Since this is a random-effects meta-analysis, this tells us that the average effect size in the universe of comparable studies is estimated as 0.506 and probably falls in the range of 0.361 to 0.650. The Z-value of 6.862 and the corresponding *p*-value of <0.001 test the null hypothesis that the average effect size in the universe of comparable studies is precisely zero. We can reject the null hypothesis and conclude that the average effect size is greater than zero – that the treatment is helpful. At the right [H], the program displays the relative weight assigned to each study when computing the combined effect size.

## HOW MUCH DOES THE EFFECT SIZE VARY?

The *average* effect size represents a substantial clinical improvement. But to understand the potential utility of this intervention, we need to also know how much the effect size varies across populations. Is the intervention consistently effective or is the impact trivial in some populations and exceptional in others? Is the intervention always beneficial, or is it sometimes harmful?

To address these questions, we can click a tool on the menu bar [I] in Figure 49.2 and display the tables shown in Figure 49.3. The statistics at the top of Figure 49.3 [J] are the same as those in Figure 49.2 and address the *average* effect size. The statistics

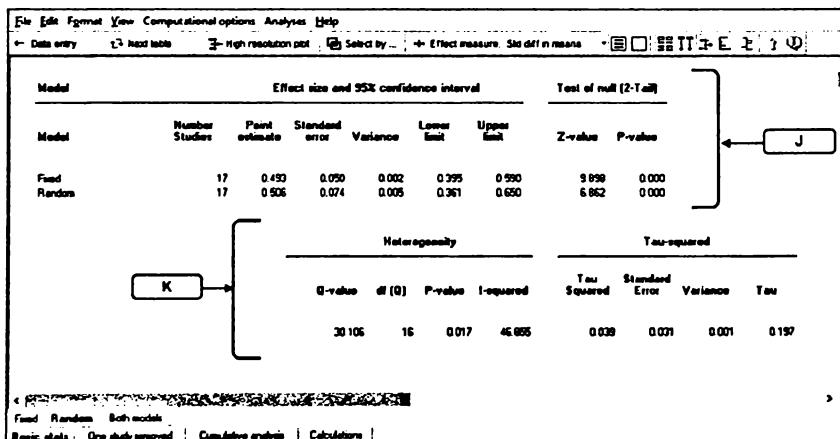


Figure 49.3 Average effect size (top), variation in effect size (bottom).

at the bottom of Figure 49.3 [K] address the *variation* in effect size across studies, as follows.

We can test the null hypothesis that all studies share a common effect size and that the variance in observed effects is due entirely to sampling error. The test statistic  $Q$  is 30.106 with 16 degrees of freedom and a corresponding  $p$ -value of 0.017. We conclude that the impact of methylphenidate is stronger in some populations than in others. However, the important question is not whether the effect size varies *at all*, but rather *how much* it varies. We turn to that now.

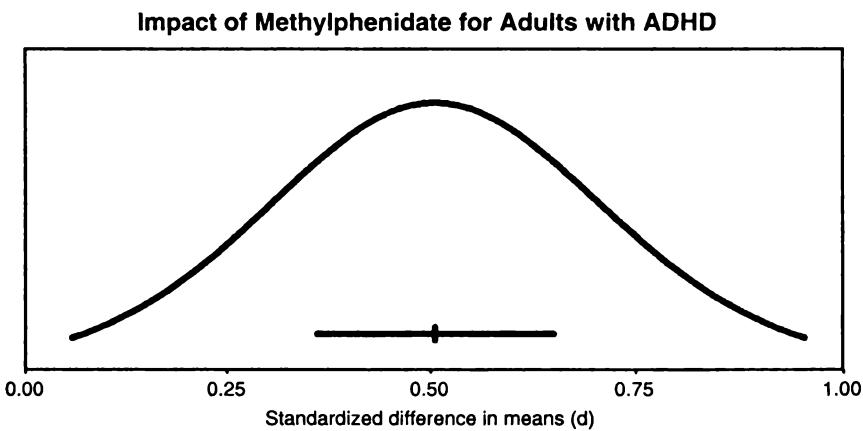
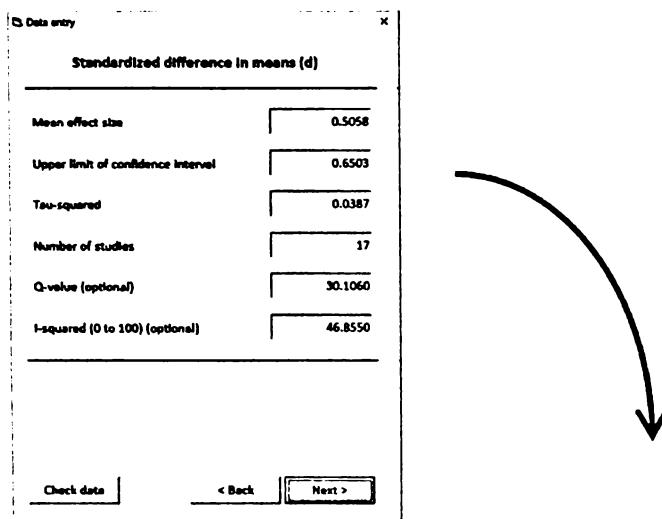
To get a general sense of the dispersion, we can start with the forest plot (Figure 49.2), where the observed effect sizes vary from -0.26 to +1.30. However, only some of this dispersion reflects variation in true effects (the variation that we care about), while the rest reflects variance due to sampling error. Returning to Figure 49.3, the  $I^2$  statistic is the ratio of  $V_{TRUE}$  to  $V_{OBSERVED}$ , and as such, it provides some context for understanding the forest plot. When  $I^2$  is low, the variance in the forest plot is mostly due to sampling error. When  $I^2$  is high, the variance in the forest plot provides a reasonable estimate for the variance of true effects. Here,  $I^2$  is around 47%, so the variation in true effects is somewhat less than the variation displayed in the forest plot. (There is a common belief that  $I^2$  tells us how much the effect size varies, but this belief is incorrect. As explained in Chapters 19 and 20,  $I^2$  is a proportion, not an absolute value.) The program displays  $T^2$ , the estimated variance of true effects (0.039) and  $T$ , the standard deviation of true effects (0.197).

## PLOT SHOWING DISTRIBUTION OF EFFECTS

While most reports of meta-analyses tend to highlight the statistics outlined above, none of these statistics directly addresses the question ‘What is the expected range of true effects for populations similar to those in the analysis?’

For that, we turn to the prediction interval. Version 4 of CMA (planned for release in 2021) will create the plot shown in Figure 49.4. The caption, *The true effect size in 95% of all populations falls in the interval 0.06 to 0.95*, is generated automatically. The plot displays the corresponding distribution of effects, which allows us to gauge the approximate proportion of effects in any given range. In this case, all effects fall to the right of zero (there are no populations where the treatment is harmful). Additionally, if we assume that 0.35 is the lower bound of a clinically important effect, we would conclude that the effect is clinically useful in more than 70% of all comparable populations.

A stand-alone program to create this plot may also be downloaded on the book’s website. This program can be used with Review Manager™ and other programs, as well as CMA. We would enter the number of studies (17), the mean effect size (0.5058), the upper limit of the mean effect size (0.6503), and  $T^2$  (0.0387). The program then generates the plot shown in Figure 49.4. Beginning with version 4 of CMA, this plot will be integrated into CMA. The stand-alone version will still be available for those using other software.



The mean effect size is 0.51 with a 95% confidence interval of 0.36 to 0.65  
 The true effect size in 95% of all comparable populations falls in the interval 0.06 to 0.95

**Figure 49.4** Plotting distribution of true effects. ADHD.

### HIGH-RESOLUTION PLOT

Click the menu button labeled 'High-resolution plot' to create the plot displayed in Figure 49.5. Menus allow the user to extensively customize the plot and then export a copy directly to Microsoft<sup>TM</sup> Word<sup>TM</sup> or PowerPoint<sup>TM</sup>.

### Methylphenidate for Adults with ADHD

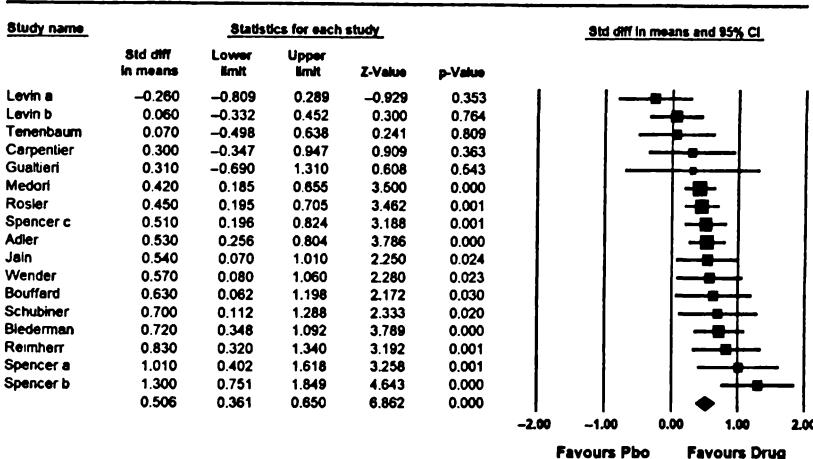


Figure 49.5 High-resolution plot in CMA.

### SUBGROUP ANALYSIS

At this point, we have established that methylphenidate is more effective in some populations than in others, and we might want to identify factors associated with the magnitude of the effect. One possible factor is the nature of the population. Specifically, some studies excluded patients who were abusing drugs, while others enrolled these patients. We want to compute the effect size separately for each subgroup of studies and then to compare the two values (see Chapter 21).

We return to the analysis screen and use the ‘Computational options’ menu to ‘Group by > SUD’ [N]. The result is displayed in Figure 49.6.

First, we assess the mean effect of treatment for each subgroup of studies. For studies that excluded drug abusers [L], the combined effect size is 0.577 with a 95% confidence interval of 0.438 to 0.717 a Z-value of 8.090 and a p-value of < 0.001. For studies that included drug abusers [M], the combined effect size is 0.162 with a 95% confidence interval of -0.136 to +0.460, a Z-value of 1.064, and a p-value of 0.287.

Next, we want to compare the effect size in the two subgroups. That is, we want to ask whether the treatment’s impact is different in studies that exclude drug abusers as compared with studies that include drug abusers. A button on the menu bar allows us to switch between the plot in Figure 49.6 and the details in Figure 49.7.

For the analyses comparing the impact of methylphenidate in studies that included drug abusers vs. studies that excluded drug abusers, we use a mixed-effects model, at the bottom of Figure 49.7. The subgroups are fixed, in the sense that we are comparing these two drugs specifically, and not generalizing to any other drugs. Within each

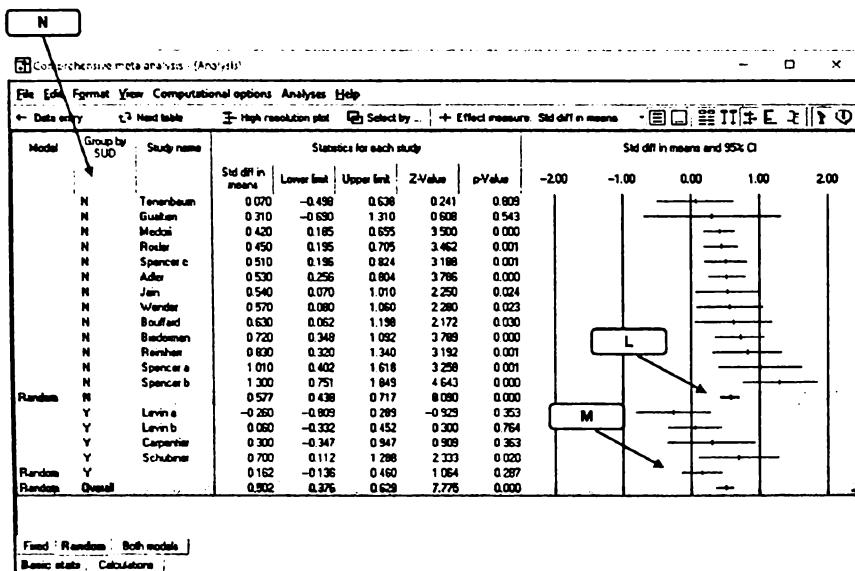


Figure 49.6 Impact of treatment as a function of subgroup: Forest plot.

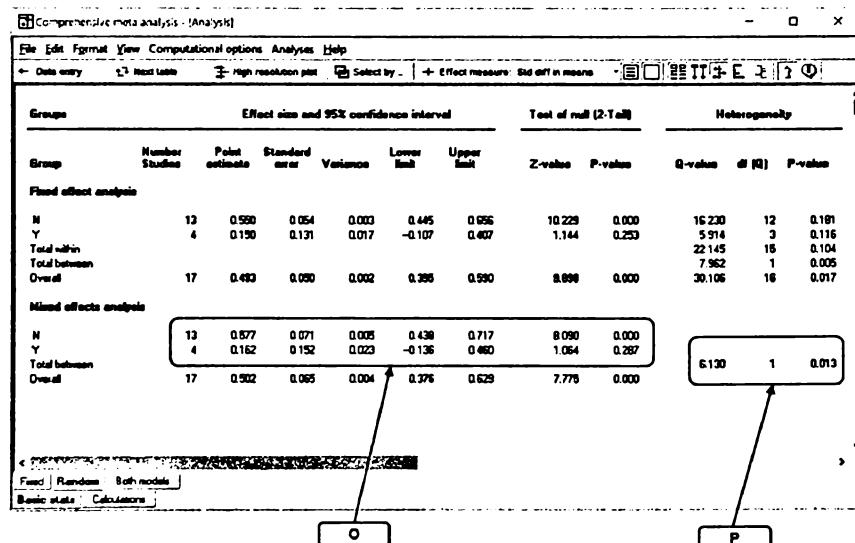


Figure 49.7 Impact of treatment as a function of subgroup: Statistics.

subgroup, the studies are a random sample of all relevant studies. Since the model is fixed at one level and random at the other level, it is called mixed-effects (Chapter 21).

The combined effect size for studies that exclude SUD patients is 0.577 with a 95% confidence interval of 0.438 to 0.717, while the combined effect size for studies that include SUD patients is 0.162 with a 95% confidence interval of -0.136 to +0.460 [O]. To test the difference between the two effect sizes, we may use a *Q*-test. The *Q*-value for this difference is 6.130 with 1 degree of freedom and a *p*-value of 0.013 [P]. We conclude that the treatment is more effective in the subgroup of populations that exclude SUD patients and less effective in the subgroup of populations that includes these patients.

It is important to recognize that (with rare exceptions) subgroup comparisons in a meta-analysis are observational by nature and cannot prove a causal relationship (Chapter 23). In this example, it is *possible* that methylphenidate is more effective in populations that exclude SUD patients because it actually works better in these patients, which *would be* a causal relationship. But it is also possible that methylphenidate was more effective in the studies which excluded SUD patients for other reasons. For example, it is possible that these studies tended to employ a higher dose of methylphenidate, and it is the higher dose (rather than the fact that these patients abused drugs) that was responsible for the larger effect in these studies. We can use meta-regression to assess the relationship between SUD and effect, with Dose held constant. We turn to that now.

## META-REGRESSION

In a primary study, we may use regression analysis to study the relationship between covariates and outcome. Similarly, in a meta-analysis we may use regression to study the relationship between covariates and effect size. In this case, the procedure is commonly called meta-regression. In a primary study, the unit of analysis is the *individual*, with covariates and outcome measured for each individual. In a meta-analysis, the unit of analysis is the study, with covariates and outcome measured for each study. However, with some modifications, the full arsenal of procedures that fall under the heading of ‘regression’ in primary studies is also available in meta-analysis (Chapter 22).

In the current example, we want to see whether the impact of methylphenidate is related to whether the study excluded patients who were abusing drugs (SUD) and/or the mean dose employed in the study (Dose). On the main analysis screen, we select ‘Meta-regression 2’ on the ‘Analyses’ menu. We define a regression with these two covariates, and the program displays the results in Figure 49.8. The results based on the random-effects model are shown here. The user may also choose to use a fixed-effect model though this is generally not recommended. The analyses displayed here are based on the Knapp–Hartung adjustment (Chapter 26).

The table at the top [Q] provides details for the relationship between SUD and effect size, with Dose held constant. The coefficient for SUD tells us that (with dose held constant) the mean effect size for studies that enrolled SUD patients is 0.4492 *d* lower than

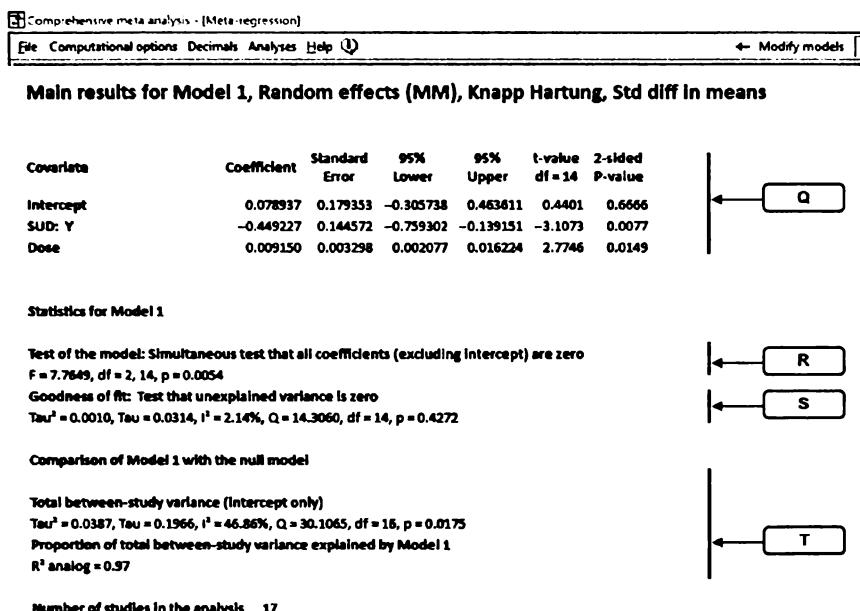


Figure 49.8 Results for regression, random effects.

for studies which excluded those patients. The confidence interval for the coefficient is  $-0.7593$  to  $-0.1392$ , and a test of the null hypothesis that SUD is not related to effect size yields a  $t$ -value of  $-3.1073$  with  $14$   $df$ , and a corresponding  $p$ -value of  $0.0077$ . This tells us that the relationship between SUD and effect size is not due to a confound with dose.

The coefficient for Dose is displayed as  $0.0092$ . This tells us that for every one-unit increase in dose, the effect size will increase by approximately  $0.01$ . The  $95\%$  confidence interval for the coefficient is  $0.0021$  to  $0.0162$ , and a test of the null hypothesis that dose is not related to effect size yields a  $t$ -value of  $2.7746$  with  $14$   $df$  and a  $p$ -value of  $0.0149$ . This is plotted in Figure 49.9, where we see that the treatment is more effective in studies that employed a higher dose of the drug. Concretely, as the dose increases from 30 units to 82 units [points U to V on the regression line], the impact of treatment increases by 48 points.

As noted, in Figure 49.8, the table at the top [Q] displays statistics for the *unique* impact of *each* covariate. By contrast, the other sections display statistics for the *joint* impact of *all* covariates.

In the section labeled 'Test of Model' [R], we test the null hypothesis that *none* of the covariates explains any variation in effect size. The  $F$ -value for this test is  $7.7649$  with  $2, 14$  degrees of freedom and a corresponding  $p$ -value of  $0.0054$ . We reject the null hypothesis and conclude that at least one of the covariates is related to effect size.

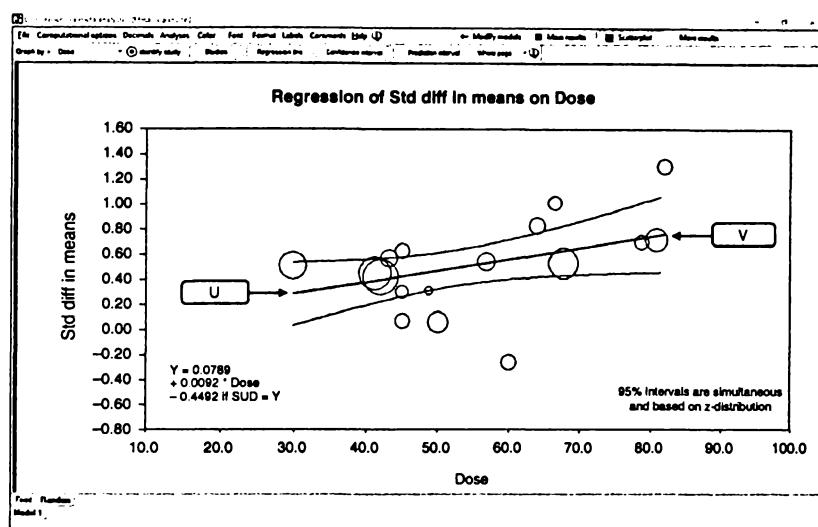


Figure 49.9 Regression of effect size on Dose, with SUD held constant.

The section labeled 'Goodness of fit' [S] addresses the residual variance. The estimated variance of true effects about the regression line ( $T^2$ ) is 0.0010, and the standard deviation of true effects about the regression line ( $T$ ) is 0.0314. The  $I^2$  statistic is 2.14%, which tells us that some 2% of the observed variance about the regression lines based on these covariates reflects variation in true effects rather than sampling error. The test for heterogeneity yields a  $Q$ -value of 14.31 with 14 degrees of freedom and a corresponding  $p$ -value of 0.4272. We conclude that the variance of observed effects about the regression line could be due entirely to sampling error. Finally, the program displays  $R^2$  as 0.97, which tells us that the covariates can explain some 97% of the initial variance in true effects.

The program offers a number of options for regression. It allows the user to include categorical covariates in the model (as it did here for SUD). In this case, the program will automatically create a set of dummy variables to represent the covariate. It allows the user to select either the Z-distribution or the Knapp–Hartung adjustment for computing confidence intervals and  $p$ -values (Knapp & Hartung, 2003) as discussed in chapter 26. It allows the user to estimate  $\tau^2$  using either the method of moments, maximum likelihood, or restricted maximum likelihood. It allows the user to define sets of covariates (for example, the linear and curvilinear impact of dose) and to assess the impact of the full set with other covariates held constant. It allows the user to define multiple prediction models and then compare them with each other.

As was true for analyses that compared subgroups, the relationships explored in meta-regression (with rare exceptions) are observational rather than causal. In this example, we attempted to identify the relationship between SUD and effect size while

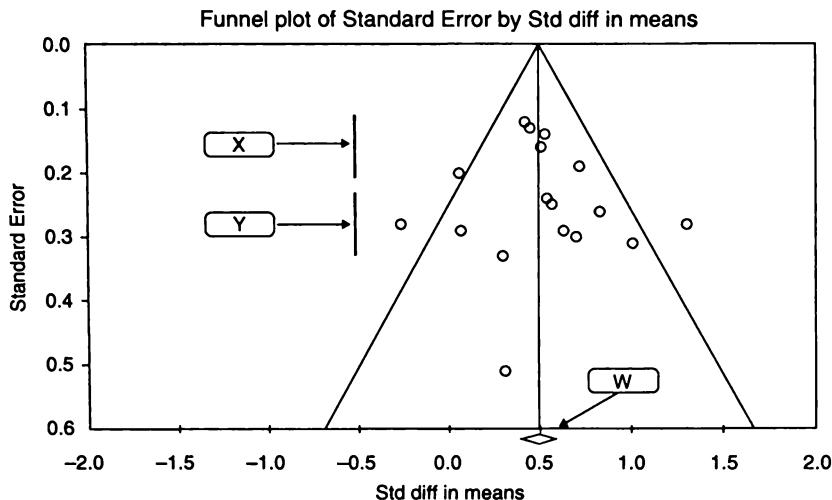


Figure 49.10 Funnel plot of observed effects.

controlling for dose and vice versa, but there may be other confounding variables that we have not considered.

## PUBLICATION BIAS

To address the potential impact of publication bias, we can select 'Analyses > Publication bias' on the main analysis screen, to display Figure 49.10. This figure shows the effect size (on the  $x$ -axis) by the standard error (on the  $y$ -axis). The large studies appear at the top, and the smaller studies appear toward the bottom.

In this case, the sample size in most studies falls within a narrow range, and so it is not likely that the procedures normally employed to assess potential for publication bias would be effective. For purposes of this volume, we can nevertheless apply the Trim and Fill procedure (Chapter 35).

A vertical line denotes the average effect size [W]. If the effects are normally distributed, we would expect half the studies to fall on either side of the line. This is the case toward the top [X] but as we move toward the smaller studies [Y] there is a predominance of studies on the right and relatively few on the left. One possible reason for this could be that the studies toward the left were not statistically significant and therefore were not published and did not find their way into the analysis. In that case, the combined effect size is based on a biased subset of all actual studies and overestimates the true average effect size.

The Trim and Fill method (Duval & Tweedie, 2000b) employs an iterative procedure to identify the studies that may be missing. It then 'creates' these studies and inserts them into the analysis. These are displayed here as filled circles [Z] which are the

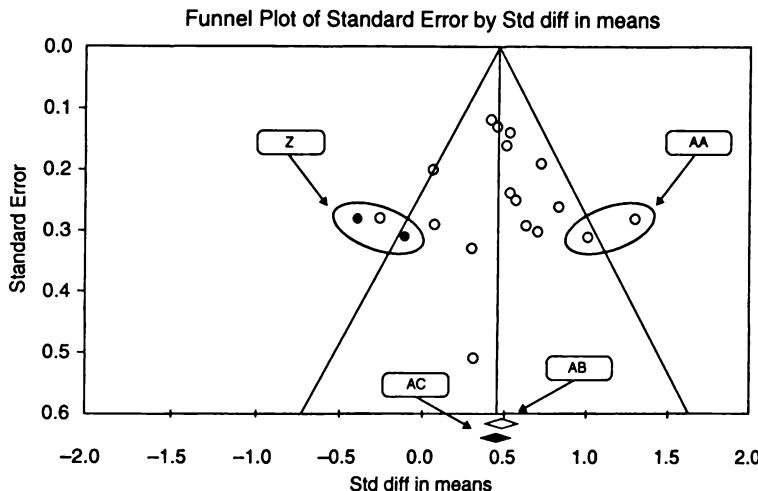


Figure 49.11 Funnel plot of observed and imputed effects.

mirror image of the actual studies [AA]. We can use all the studies (actual and imputed) to compute an adjusted estimate of the mean effect size. The initial estimate of the combined effect size was 0.506 [AB] but the adjusted value (included the imputed studies) is 0.442 [AC]. If the asymmetry was due to publication bias, then this adjustment yields an estimate of the unbiased effect size. The adjustment is minor, in the sense that the clinical meaning of 0.442 is essentially the same as the clinical meaning of 0.506. Note that there are reasons other than publication bias that may explain or contribute to asymmetry in funnel plots. As explained in Chapter 35, the Trim and Fill procedure should be seen as a sensitivity analysis. It is not intended to yield a "correct" effect size.

CMA also features other methods that are typically used with the aim of testing and/or adjusting for publication bias. These include the Egger test of the intercept, the Begg and Mazumdar rank correlation test, and Rosenthal's Fail-safe  $N$  (Rothstein *et al.*, 2005; Sterne *et al.*, 2011). The program can also generate a text report that explains how to interpret the results for each of the publication bias procedures.

## EXPLAINING RESULTS

The following is how one might explain the results of this analysis.

### Overview

This example is a re-analysis of a systematic review published by Castells *et al.* (2011). The analysis is based on seventeen studies that evaluated the effect of methylphenidate on cognitive function in adults with attention deficit hyperactivity disorder (ADHD).

In each study, patients were randomly assigned to either drug or a placebo and the researchers assessed the patients' cognitive function at the conclusion of treatment. The effect size index is the standardized mean difference ( $d$ ). The results of this analysis will be generalized to comparable studies, and so the random-effects model was employed for the analysis.

In this context, a standardized mean difference of 0.20 would be considered trivial – this is a difference that shows up on the tests, but the patient might not be aware of any change. A standardized mean difference of 0.50 would be considered moderate – the patient would recognize that they were doing better than usual, and coworkers might be aware of a change. A standardized mean difference of 0.80 would be considered large – the patient would feel great, and the difference would be obvious enough that others might remark on it.

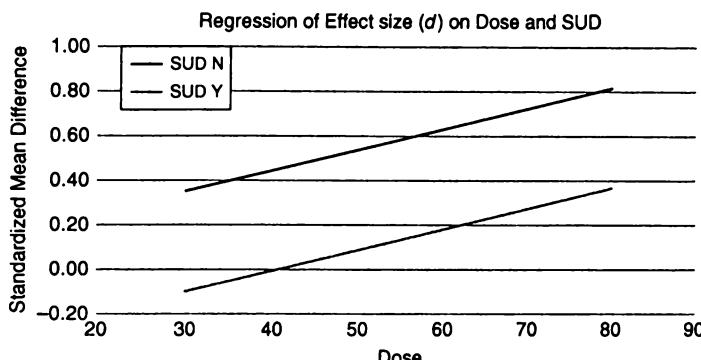
### Does methylphenidate affect cognitive scores?

The standardized mean difference is 0.506. *On average*, methylphenidate increased cognitive functioning by 0.506 standard deviations as compared with placebo. The confidence interval for the standardized mean difference is 0.361 to 0.650, which tells us that the mean effect size in the universe of comparable studies could fall anywhere in this range. This range does not include an effect size of zero, which tells us that the mean effect size is probably not zero. Similarly, the  $Z$ -value for testing the null hypothesis (that  $d$  is 0.0) is 6.862, with a corresponding  $p$ -value of  $<0.001$ . Using the Knapp–Hartung adjustment,  $t = 6.29$ ,  $df = 16$ ,  $p < 0.001$ , and the 95% confidence interval is 0.335 to 0.676 (see Chapter 26). We can reject the null hypothesis and conclude that (on average) the drug does increase cognitive function in the universe of populations which are comparable to those in the analysis. Given the dispersion in effects (as discussed below), it is important to recognize that the mean effect size applies to this particular mix of studies and would be different for another mix of populations, dosages, and so on.

### How much does the effect size vary across studies?

The  $Q$ -statistic provides a test of the null hypothesis that all studies in the analysis share a common effect size. If all studies shared the same effect size, the expected value of  $Q$  would be equal to the degrees of freedom (the number of studies minus 1). The  $Q$ -value is 30.106 with 16 degrees of freedom and  $p=0.017$ . We reject the null hypothesis that the true effect size is identical in all the studies. The  $I^2$  statistic is 47%, which tells us that 47% of the variance in observed effects reflects variance in true effects rather than sampling error. The variance of true effects ( $T^2$ ) is 0.039, and the standard deviation of true effects ( $T$ ) is 0.197.

The prediction interval is 0.058 to 0.953. We would expect that in some 95% of all populations comparable to those in the analysis, the true effect size will fall in this range. Based on the context outlined above, there will be some populations where the impact of the treatment is trivial and some where it is large.



**Figure 49.12** Regression of effect size ( $d$ ) on Dose and SUD. Plot created in Excel (TM).

### Is effect size related to Dose and/or SUD?

For every one mg. increase in dose there is an increase of roughly one point in the effect size,  $d$ . As dose increases by 50 mg. (from 0.30 to 0.80) the predicted effect size increases by 46 points. This is reflected by the slope of the prediction lines in Figure 49.12. For studies that enrolled SUD patients, the predicted effect size increases from  $-0.10$  to  $+0.36$ . For studies that excluded SUD patients, the predicted effect size increases from  $0.35$  to  $0.81$ .

The predicted effect size for studies that enrolled SUD patients is 45 points lower than for studies that excluded these patients. This is reflected in the difference between the two prediction lines in Figure. The bottom line (for studies that enrolled SUD patients) is 45 lower than the top line (for studies that included these patients).

The relationship between dose and effect size remains even after we partial SUD, and the relationship between SUD and effect size remains even after we partial dose. The coefficient for Dose as a predictor of effect size (with SUD held constant) is 0.0092 with a 95% confidence interval of 0.0021 to 0.1162. For a test of the null hypothesis that there is no relationship between dose and effect size,  $t(14) = 2.7746$  and  $p = 0.0149$ . The coefficient for SUD (present) as a predictor of effect size (with dose held constant) is  $-0.4492$  with a 95% confidence interval of  $-0.7593$  to  $-0.1392$ . For a test of the null hypothesis that there is no relationship between dose and effect size,  $t(14) = -3.1073$  and  $p = 0.0077$ .

### Publication bias

While it is likely that some studies are missing from the analysis due to publication bias (since this is typically the case), the impact of missing studies in this analysis was probably minor. The Trim and Fill analysis suggests that there may be two missing studies, but if we impute these studies and include them in the analysis, the mean effect size shifts only slightly (from 0.506 to 0.442).



# How to Explain the Results of an Analysis

---

Introduction

The overview

The mean effect size

Variation in effect size

Notations

Impact of resistance exercise on pain

Correlation between letter knowledge and word recognition

Statins for prevention of cardiovascular events

Bupropion for smoking cessation

Mortality following mitral-valve procedures in elderly patients

---

## INTRODUCTION

In this chapter, we provide examples of how one might explain the results of a simple meta-analysis, for example to a colleague. There is one example based on each of several effect sizes (a standardized difference in means, a risk ratio, an odds ratio, a correlation, and the risk in one group). The explanations might form the basis for a formal write-up of the results in, say, a journal article.

In each case, we show how to take the data reported by a computer program and use this to write a report. The examples shown are based on the software comprehensive meta-analysis (CMA) which was developed by the authors of this text. However, any software will report the same statistics.

In each example, there are sections labeled overview, mean effect size, and heterogeneity. Here, we do not address other issues such as publication bias.

The analyses shown here are selected to illustrate various statistical issues and should not be used to provide information about the efficacy or safety of any intervention. Any reader interested in these interventions should study the original papers.

## THE OVERVIEW

In this section, we introduce the analysis by providing some basic information such as the number of studies and the effect-size index. Here, we also provide the rationale for using the random-effects model. The decision to use this model is based on the fact that we will be making an inference to a universe of comparable studies and is not informed by any statistical test (see Chapter 13).

## THE MEAN EFFECT SIZE

We present the mean effect size, the confidence interval for the mean, and a test of the null hypothesis that the mean effect size is zero (for a difference) or one (for a ratio).

## VARIATION IN EFFECT SIZE

When researchers ask about heterogeneity, they intend to ask, ‘How much does the effect size vary across studies?’. Unfortunately, the statistics typically reported for heterogeneity ( $Q$ ,  $p$ ,  $I^2$ ,  $T^2$ ) do not directly address this question. We may need to report these statistics (since editors expect to see them), but we can also explain what they mean, so that they are interpreted correctly. The one statistic that does address the question of interest is the prediction interval. We report this interval, along with a plot that displays the entire distribution of true effects. We then show how this informs the discussion of the potential utility of the intervention.

The analyses that follow each have at least fourteen studies, and therefore provide a useful estimate of the prediction interval. That interval is indicated on each forest plot by the line labeled [H], and also by a separate plot that displays the distribution of true effects. When there are only a few studies in the analysis, it may not be possible to compute a reliable estimate of the prediction interval. In that case, we would either omit the interval or explain that the numbers are not reliable.

## NOTATIONS

Within each report, each statistic is followed by a notation such as [A], [B], and so on. These refer to the screenshots that follow the report and show where one would find these statistics in CMA. The screenshots do not show the computation of the Knapp–Hartung adjustment, but this is discussed in chapter 26. The plots showing the distribution of effects were created by a program that is available for download on the book’s website ([www.Introduction-to-Meta-Analysis.com](http://www.Introduction-to-Meta-Analysis.com)).

## Datasets

The datasets for these examples, along with additional details of each analysis, can be downloaded at the book’s website.

## IMPACT OF RESISTANCE EXERCISE ON PAIN

This example is based on data in a systematic review published by Juhl, Christensen, Roos, Zhang, and Lund (2014). References in the text refer to Figure 50.1 through Figure 50.4.

### Overview

The analysis is based on 32 studies that evaluated the effect of resistance exercise on pain from knee osteoarthritis. In each study, patients were randomly assigned to either the exercise condition or to a control condition, and patients rated their degree of pain after the treatment.

The effect-size index is the standardized mean difference ( $g$ ). A positive score reflects a reduction in pain. In this context, a standardized mean difference of 0.20 would be considered small, a standardized mean difference of 0.50 would be considered moderate, and a standardized mean difference of 0.80 would be considered large. Clinicians assess pain using the visual analog scale (VAS) where patients rated their pain on a line with endpoints of zero and one hundred. In a population where the standard deviation of the VAS was 20 points, these three effect sizes would correspond to difference in VAS of 4 points, 10 points, and 16 points, respectively.

The results of this analysis will be generalized to comparable studies, and so the random-effects model was employed for the analysis.

### Does resistance exercise affect the pain reported by patients?

The standardized mean difference is 0.629 [A]. *On average*, the exercise reduced pain 0.629 standard deviations as compared with placebo. The confidence interval for the standardized mean difference is 0.458 to 0.799 [B], which tells us that the mean effect size in the universe of comparable studies could fall anywhere in this range. This range does not include an effect size of zero, which tells us that the mean effect size is probably not zero. Similarly, the Z-value for testing the null hypothesis (that  $g$  is 0.0) is 7.235, with a corresponding  $p$ -value of  $<0.001$  [C]. Using the Knapp–Hartung adjustment,  $t = 5.654$ ,  $df = 31$ ,  $p < 0.001$ , and the 95% confidence interval is 0.402 to 0.856. We can reject the null hypothesis and conclude that (on average) resistance exercise does decrease pain in the universe of populations which are comparable to those in the analysis. Given the dispersion in effects (as discussed below), it is important to recognize that the mean effect size applies to this particular mix of studies and would be different for another mix of populations or variants of the intervention.

### How much does the effect size vary across studies?

The  $Q$ -statistic provides a test of the null hypothesis that all studies in the analysis share a common effect size. If all studies shared the same effect size, the expected value of  $Q$  would be equal to the degrees of freedom (the number of studies minus 1).

The  $Q$ -value is 103.110 with 31 degrees of freedom and  $p < 0.001$ . We reject the null hypothesis that the true effect size is identical in all the studies [D]. The  $I^2$  statistic is 69.935% [E], which tells us that roughly 70% of the variance in observed effects reflects variance in true effects rather than sampling error. The variance of true effects ( $T^2$ ) is 0.151 [F], and the standard deviation of true effects ( $T$ ) is 0.389 [G].

The prediction interval is  $-0.185$  to  $+1.442$  [H]. We would expect that in some 95% of all populations comparable to those in the analysis, the true effect size will fall in this range. The distribution of effects (Figure 50.4) shows that the effect will be positive (greater than zero) in roughly 94% of populations. Based on the context outlined above, the impact will be moderate ( $g = 0.50$ ) or higher in roughly 61% of all populations, and large ( $g = 0.80$ ) or higher in roughly 32% of all populations. If the effects are normally distributed, we would expect the intervention to increase pain in some 6% of all comparable populations.

Run analyses →												
	Study name	Group-A Mean	Group-A Std-Dev	Group-A Sample size	Group-B Mean	Group-B Std-Dev	Group-B Sample size	Effect direction	Hedges's g	Std Err	Variance	
1	Cheung	37.800	64.000	15	49.600	42.400	16	Auto	-0.213	0.351	0.123	
2	Weidenbaum	0.100	2.300	19	0.100	1.000	20	Auto	0.000	0.314	0.098	
3	Sayers b	1.800	3.400	10	1.500	2.600	6	Auto	0.090	0.488	0.239	
4	Sayers a	1.800	2.800	12	1.500	2.600	6	Auto	0.104	0.477	0.227	
5	McKnight	1.350	9.000	95	0.000	9.000	87	Auto	0.149	0.148	0.022	
6	Meurer	43.540	86.950	49	28.490	86.950	49	Auto	0.172	0.201	0.040	
7	Berjesson	0.400	2.000	34	0.000	1.400	34	Auto	0.229	0.241	0.058	
8	Swank	0.900	7.300	36	-0.800	7.300	35	Auto	0.230	0.236	0.056	
9	Rooks	0.100	2.300	14	-0.700	4.000	15	Auto	0.236	0.363	0.131	
10	Ettlinger	0.190	0.680	120	0.000	0.560	64	Auto	0.302	0.155	0.024	
11	Frougghi	1.870	1.730	18	1.200	1.730	19	Auto	0.379	0.325	0.106	
12	McCarthy	2.130	2.910	104	0.950	2.910	86	Auto	0.404	0.147	0.022	
13	Huang b	1.200	1.600	30	0.500	1.700	32	Auto	0.418	0.254	0.064	
14	Lim	9.000	16.620	50	1.810	16.800	47	Auto	0.427	0.204	0.042	
15	Petrella	0.510	0.150	91	0.440	0.150	88	Auto	0.465	0.151	0.023	
16	Topp	1.540	3.210	57	-0.020	3.190	35	Auto	0.483	0.210	0.044	
17	Ergenekon	0.460	1.600	18	-0.500	2.200	20	Auto	0.484	0.323	0.104	
18	Weng	1.100	1.600	31	0.100	1.500	33	Auto	0.638	0.253	0.064	
19	Jan a	3.700	3.640	34	1.200	3.640	15	Auto	0.676	0.313	0.098	
20	Jan b	3.000	2.620	34	1.200	2.620	15	Auto	0.676	0.313	0.098	
21	Bezalel	3.000	4.130	25	0.000	4.130	25	Auto	0.715	0.287	0.083	
22	Baker	79.000	87.960	22	20.000	71.050	22	Auto	0.725	0.306	0.094	
23	Bennell	2.600	2.440	39	0.480	2.430	37	Auto	0.862	0.238	0.056	
24	Lin	4.600	3.310	36	1.200	4.000	18	Auto	0.944	0.239	0.089	
25	Regind	3.000	5.240	11	-1.000	2.620	12	Auto	0.944	0.426	0.181	
26	Huang a	1.600	1.500	91	0.200	1.300	33	Auto	0.959	0.211	0.044	
27	Cheng	2.300	1.300	24	0.900	1.500	17	Auto	0.991	0.330	0.109	
28	Horstmann	1.790	1.200	19	0.530	0.700	19	Auto	1.256	0.349	0.122	
29	Schäke	6.100	4.030	10	-0.400	4.030	10	Auto	1.545	0.493	0.243	
30	Safi b	3.600	1.400	24	0.600	1.300	12	Auto	2.144	0.428	0.183	
31	Safi a	4.300	1.200	23	0.600	1.300	12	Auto	2.929	0.494	0.244	
32	Gü	13.700	5.000	17	-2.500	2.600	6	Auto	3.436	0.683	0.466	
33												

Figure 50.1 Impact of resistance exercise on pain. Data-entry screen.

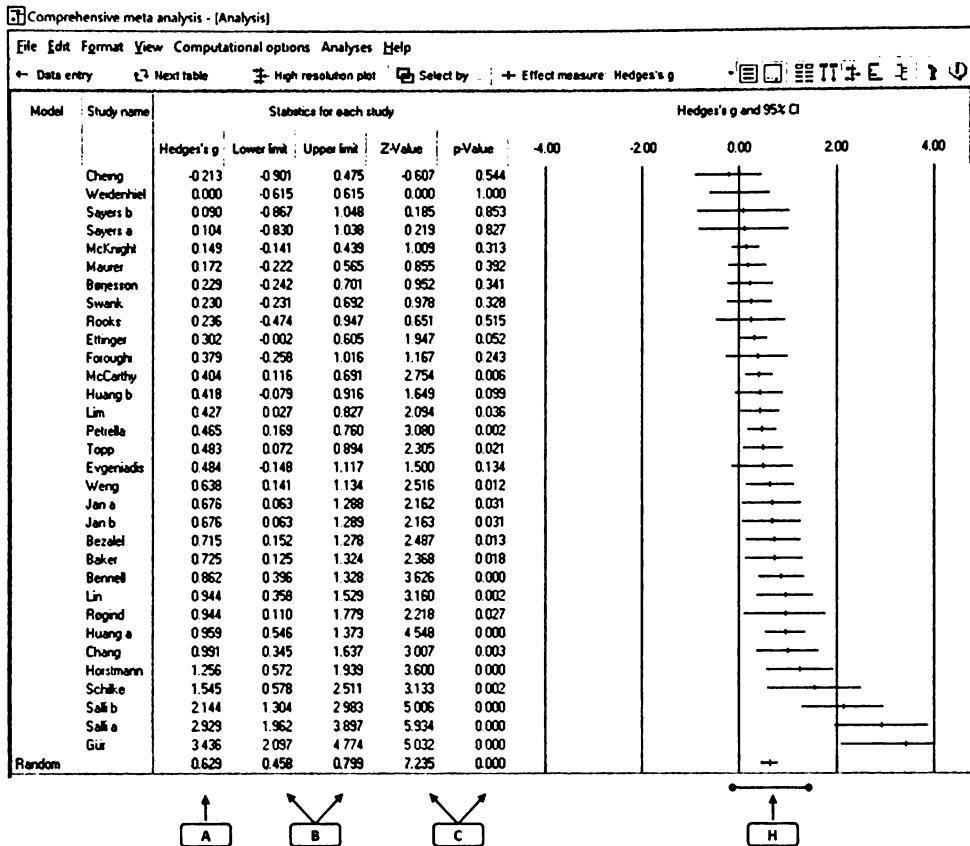
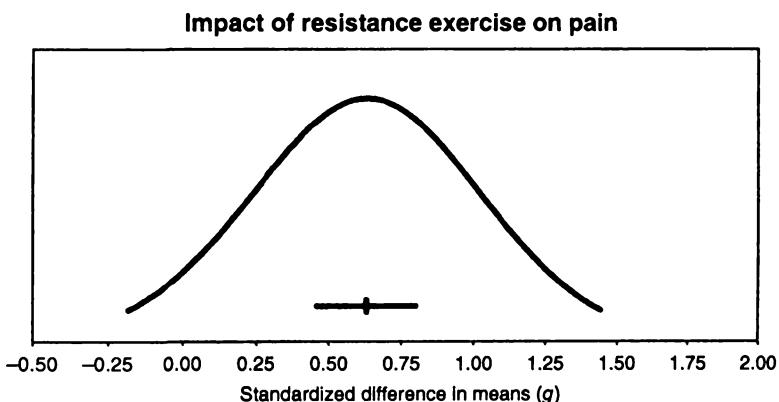


Figure 50.2 Impact of resistance exercise on pain.  $g > 0$  indicates exercise reduced pain.

Heterogeneity				Tau-squared			
Q-value	df [Q]	p-value	I-squared	Tau squared	Standard error	Variance	Tau
103.110	31	0.000	69.935	0.151	0.062	0.004	0.389

↓  
D      E      F      G

Figure 50.3 Impact of resistance exercise on pain. Heterogeneity statistics.



The mean effect size is 0.63 with a 95% confidence interval of 0.46 to 0.80  
The true effect size in 95% of all comparable populations falls in the interval -0.18 to 1.44

Figure 50.4 Impact of resistance exercise on pain. Distribution of true effects.

## CORRELATION BETWEEN LETTER KNOWLEDGE AND WORD RECOGNITION

This example is a reanalysis of a systematic review published by Hjetland, Brinchmann, Scherer, and Melby-Lervåg (2017). References in the text refer to Figure 50.5 through Figure 50.8.

### Overview

The analysis included 16 studies, each of which looked at the correlation between a child's ability to identify letters (as a preschooler) and that child's ability to recognize words (when entering school). The effect size index is the Pearson correlation coefficient,  $r$ . The results of this analysis will be generalized to comparable studies, and so the random-effects model was employed for the analysis.

For purposes of this discussion, we will assume that correlations of 0.10, 0.30, and 0.50 represent small, moderate, and large values.

### Is letter knowledge correlated with word recognition?

The mean correlation is 0.384 [A], which means that *on average*, children with higher levels of letter knowledge (as a preschooler) tended to have a higher level of word recognition (when entering school). The confidence interval for the correlation is 0.308 to 0.454 [B], which tells us that the mean correlation in the universe of comparable studies could fall anywhere in this range. This range does not include a correlation of zero, which tells us that the mean correlation is probably not zero. Similarly, the Z-value for testing the null hypothesis (that the mean correlation is zero) is 9.236, with a corresponding  $p$ -value of  $< 0.001$  [C]. Using the Knapp–Hartung adjustment,  $t = 7.980$ ,  $df = 15$ ,  $p < 0.001$ , and the 95% confidence interval is 0.288 to 0.472. We can reject the null hypothesis and conclude that in the universe of populations comparable to those in the analysis, there is (on average) a positive correlation between the two measures. Given the dispersion in correlations (as discussed below), it is important to recognize that the mean correlation applies to this particular mix of studies and would be different for another mix of populations.

### How much does the correlation vary across studies?

The  $Q$ -statistic provides a test of the null hypothesis that all studies in the analysis share a common correlation. If all studies shared the same correlation, the expected value of  $Q$  would be equal to the degrees of freedom (the number of studies minus 1). The  $Q$ -value is 63.059 with 15 degrees of freedom and a corresponding  $p$ -value of  $< 0.001$  [D]. We reject the null hypothesis that the true correlation is identical in all the studies. The  $I^2$  statistic is 76.213%, which tells us that some 76% of the variance in observed correlations reflects variance in true correlations rather than sampling error [E]. The variance of true correlations ( $T^2$ ) is 0.022 in Fisher's Z units [F], and the standard deviation of true correlations ( $T$ ) is 0.149 in Fisher's Z units [G].

The 95% prediction interval is +0.071 to +0.628 [H]. The distribution of correlations is plotted in Figure 50.8. The true correlation between the early test and later performance varies from one population to the next, and in any single population, it likely falls in this range. Using the criteria outlined above in conjunction with Figure 50.8, the correlation would be greater than zero in all populations; moderate or greater in 71% of all populations; and large or greater in 17% of all populations.

Note that the distribution is assumed to be symmetric in Fisher's Z units. It appears to be skewed because the plot uses correlations rather than Fisher's Z units on the *x*-axis.

Comprehensive meta-analysis - [C:\Users\Björn\Dropbox\S 000 Second Edition\Data Sets\000 Correlations Reading predictions Campbell\Reading.cma]

---

File Edit Format View Insert Identify Tools Computational options Analyses Help

---

Run analyses →

---

	Study name	Correlation	Sample size	Effect direction	Correlation	Std Err	Variance	Fisher's Z	Std Err	Variance	
1.	Nastlund & Schneider, 1996	-0.040	89	Auto	-0.040	0.108	0.012	-0.040	0.108	0.012	
2.	Lepola et al., 2016	0.060	90	Auto	0.060	0.107	0.011	0.060	0.107	0.011	
3.	Fricke et al., 2016	0.180	78	Auto	0.180	0.112	0.012	0.182	0.115	0.013	
4.	Sears & Keogh, 1993	0.220	104	Auto	0.220	0.095	0.009	0.224	0.100	0.010	
5.	Bishop, & League, 2006	0.290	79	Auto	0.290	0.105	0.011	0.299	0.115	0.013	
6.	Aarnoutse et al., 2005	0.340	78	Auto	0.340	0.102	0.010	0.354	0.115	0.013	
7.	Shalit & Share, 2003	0.360	313	Auto	0.360	0.049	0.002	0.377	0.057	0.003	
8.	Hecht et al., 2000	0.400	197	Auto	0.400	0.060	0.004	0.424	0.072	0.005	
9.	Schatschneider et al., 2004	0.440	189	Auto	0.440	0.059	0.003	0.472	0.073	0.005	
10.	Leppänen et al., 2008	0.450	158	Auto	0.450	0.064	0.004	0.485	0.080	0.006	
11.	Piasta et al., 2012	0.450	371	Auto	0.450	0.042	0.002	0.485	0.052	0.003	
12.	Sawyer, 1992	0.470	300	Auto	0.470	0.045	0.002	0.510	0.058	0.003	
13.	Parila et al., 2004	0.500	95	Auto	0.500	0.078	0.006	0.549	0.104	0.011	
14.	Prochnow et al., 2013	0.530	76	Auto	0.530	0.084	0.007	0.590	0.117	0.014	
15.	Bowey, 1995	0.580	116	Auto	0.580	0.062	0.004	0.662	0.094	0.009	
16.	Muter et al., 2004	0.620	90	Auto	0.620	0.066	0.004	0.725	0.107	0.011	
17.											

Figure 50.5 Predicting reading scores. Data-entry screen.

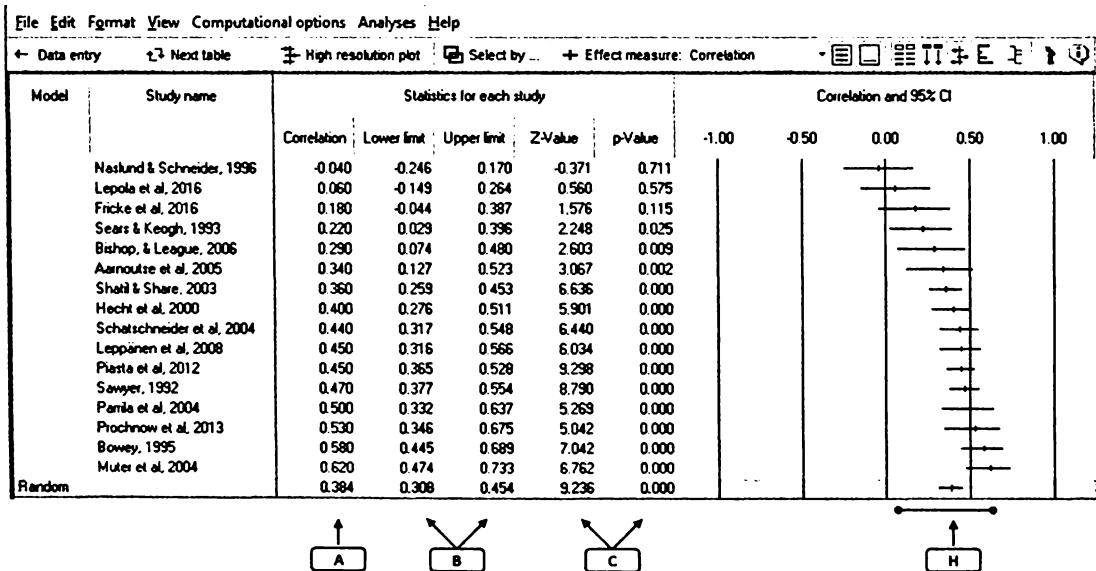
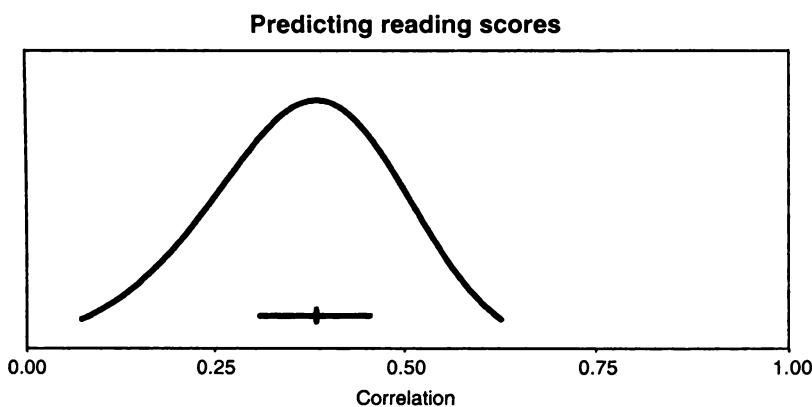


Figure 50.6 Predicting reading scores.

Heterogeneity				Tau-squared			
Q-value	df [Q]	p-value	I-squared	Tau squared	Standard error	Variance	Tau
63.059	15	0.000	76.213	0.022	0.012	0.000	0.149

D → Q-value  
 E → I-squared  
 F → Tau squared  
 G → Tau

Figure 50.7 Predicting reading scores. Heterogeneity statistics.



The mean effect size is 0.38 with a 95% confidence interval of 0.31 to 0.45  
 The true effect size in 95% of all comparable populations falls in the interval 0.07 to 0.63

Figure 50.8 Predicting reading scores. Distribution of true correlations.

## STATINS FOR PREVENTION OF CARDIOVASCULAR EVENTS

This example is a reanalysis of a systematic review published by Taylor *et al.* (2013). The original analysis used risk ratios, but we use odds ratios in this example. References in the text refer to Figure 50.9 through Figure 50.12.

### Overview

This analysis includes 14 studies where patients were randomized to treatment with either statins or a placebo. The outcome was a cardiovascular event. The effect size index is the odds ratio for statins vs. placebo, with odds ratios < 1 favoring statins. The results of this analysis will be generalized to comparable studies, and so the random-effects model was employed for the analysis.

#### Do statins reduce the odds that patients will have a cardiovascular event?

The mean effect size is an odds ratio of 0.710 [A]. On average, the odds of an event for those treated with statins was 0.710 as high as for those treated with a placebo. The confidence interval for the effect size is 0.641 to 0.786 [B], which tells us that the *mean* effect size in the universe of comparable studies could fall anywhere in this range. The Z-value tests the null hypothesis that the mean odds ratio is 1.0. The Z-value is  $-6.580$  with a corresponding *p*-value of  $< 0.001$  [C]. Using the Knapp–Hartung adjustment,  $t = -6.564$ ,  $df = 13$ ,  $p < 0.001$ , and the 95% confidence interval is 0.635 to 0.795. We reject the null hypothesis and conclude that in the universe of populations comparable to those in the analysis, the drug (on average) does reduce the odds of an event.

#### How much does the effect size vary across studies?

The *Q*-statistic provides a test of the null hypothesis that all studies in the analysis share a common effect size. If all studies shared the same effect size, the expected value of *Q* would be equal to the degrees of freedom (the number of studies minus 1). The *Q*-value is 13.455 with 13 degrees of freedom and *p* = 0.413 [D]. We cannot reject the null hypothesis that the true effect size is the same in all these studies. The *I*<sup>2</sup> statistic is 3.382%, which tells us that some 3% of the variance in observed effects reflects variance in true effects rather than sampling error [E]. *T*<sup>2</sup>, the variance of true effect sizes, is 0.001 in log units [F]. *T*, the standard deviation of true effect sizes, is 0.038 in log units [G].

If we assume that the effects are normally distributed (in log units) we can estimate that the prediction interval as 0.618 to 0.817 [H]. The true effect size for any single population will usually fall in this range. The distribution of effects is plotted in Figure 50.12. Since the prediction interval is only trivially wider than the confidence interval, we can assume that the true effect is basically the same for all comparable

populations. The mean effect size could fall anywhere within the confidence interval, but wherever it falls, all populations fall within a few points of that value.

Note that the distribution is assumed to be symmetric in log units. It appears to be skewed because the plot uses the odds ratio rather than the log odds ratio on the *x*-axis.

Comprehensive meta analysis - [C:\Users\Biost\Dropbox\\$ 000 Second Edition\Data Sets\000 OR Statins\Statins.cma]

---

File Edit Format View Insert Identify Tools Computational options Analyses Help

Run analyses →

---

	Study name	Statins Events	Statins Total N	Control Events	Control Total N	Odds ratio	Log odds ratio	Std Err	Variance	J
1	CERDIA 2004	0	103	4	79	0.001	-2.513	1.498	2.245	
2	PHYLIS 2004	1	253	3	254	0.332	-1.103	1.158	1.341	
3	JUPITER 2008	47	8901	95	8901	0.492	-0.709	0.179	0.032	
4	ACAPS 1994	5	460	9	459	0.549	-0.599	0.562	0.316	
5	PREVENT IT 2004	9	433	16	431	0.551	-0.597	0.422	0.178	
6	CARDS 2008	50	1429	74	1412	0.656	-0.422	0.187	0.035	
7	HYRIM 2007	6	142	9	143	0.657	-0.420	0.541	0.293	
8	Adult Japanese MEGA	66	3666	101	3966	0.665	-0.409	0.160	0.026	
9	WDSCOPS	390	3302	509	3293	0.733	-0.311	0.072	0.005	
10	AFCAPS/TexCAPS 1998	163	3304	215	3301	0.745	-0.295	0.107	0.011	
11	ASPEN 2006	72	959	75	946	0.943	-0.059	0.172	0.029	
12	KAPS 1995	2	214	2	212	0.991	-0.009	1.005	1.009	
13	CAIUS 1996	3	151	2	154	1.541	0.432	0.920	0.847	
14	METEOR 2010	6	700	0	281	5.269	1.662	1.469	2.159	
15										

Figure 50.9 Statins for prevention of cardiovascular events. Data-entry screen.

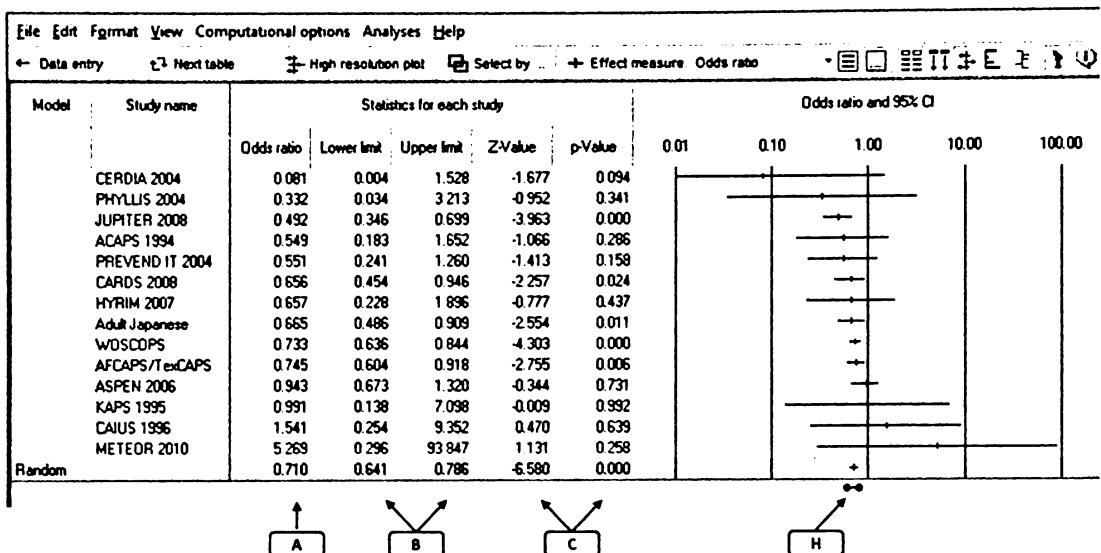
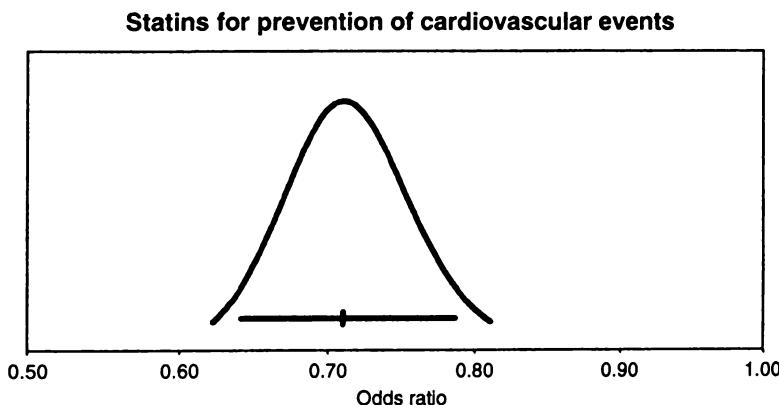


Figure 50.10 Statins for prevention of cardiovascular events. Odds ratio < 1 shows reduction in events.

Heterogeneity				Tau-squared			
Q-value	df [Q]	p-value	I-squared	Tau squared	Standard error	Variance	Tau
13.455	13	0.413	3.382	0.001	0.017	0.000	0.038

↑  
 D      E      F      G

Figure 50.11 Statins for prevention of cardiovascular events. Heterogeneity statistics.



The mean effect size is 0.71 with a 95% confidence interval of 0.64 to 0.79  
 The true effect size in 95% of all comparable populations falls in the interval 0.62 to 0.81

Figure 50.12 Statins for prevention of cardiovascular events. Distribution of true effects.

## BUPROPION FOR SMOKING CESSATION

This example is a reanalysis of a systematic review published by Hughes, Stead, Hartmann-Boyce, Cahill, and Lancaster (2014). References in the text refer to Figure 50.13 through Figure 50.16.

### Overview

This analysis includes 27 studies where patients who smoked were randomized to treatment with either Bupropion or a placebo. The outcome was a finding that the patient had stopped smoking at some point subsequent to the study. The effect size index is the risk ratio, i.e., the ratio of success for Bupropion vs. placebo. The results of this analysis will be generalized to comparable studies, and so the random-effects model was employed for the analysis.

### Does Bupropion increase the likelihood that patients will stop smoking?

The mean effect size is a risk ratio of 1.562 [A]. On average, patients treated with Bupropion were around 50% more likely to report success as compared with patients treated with a placebo. The confidence interval for the effect size is 1.364 to 1.790 [B], which tells us that the mean effect size in the universe of comparable studies could fall anywhere in this range. The Z-value tests the null hypothesis that the mean risk ratio is 1.0. The Z-value is 6.435 with a corresponding *p*-value of < 0.001 [C]. Using the Knapp–Hartung adjustment,  $t = 6.247$ ,  $df = 25$ ,  $p < 0.001$ , and the 95% confidence interval is 1.349 to 1.809. We reject the null hypothesis and conclude that in the universe of populations comparable to those in the analysis, the drug (on average) does increase the likelihood of success. Given the dispersion in effects (as discussed below), it is important to recognize that the mean effect size applies to this particular mix of studies and would be different for another mix of populations, duration of the intervention, and so on.

### How much does the effect size vary across studies?

The *Q*-statistic provides a test of the null hypothesis that all studies in the analysis share a common effect size. If all studies shared the same effect size, the expected value of *Q* would be equal to the degrees of freedom (the number of studies minus 1). The *Q*-value is 42.011 with 26 degrees of freedom and  $p = 0.024$  [D]. We can reject the null hypothesis that the true effect size is the same in all these studies. The  $I^2$  statistic is 38.112%, which tells us that some 38% of the variance in observed effects reflects variance in true effects rather than sampling error [E].  $T^2$ , the variance of true effect sizes, is 0.045 in log units [F].  $T$ , the standard deviation of true effect sizes, is 0.213 in log units [G].

If we assume that the effects are normally distributed (in log units), we can estimate that the prediction interval as 0.986 to 2.477 [H]. The distribution of effects is shown

in Figure 50.16. The true effect size for any single population will usually fall in this range. This means that the treatment will increase the likelihood of success by some amount in almost all comparable populations; by at least 50% in 55% of all comparable populations; and by at least 100% in 13% of all comparable populations.

Note that the distribution of effects is assumed to be symmetric in log units. It appears to be skewed because the plot uses the risk ratio rather than the log risk ratio on the  $x$ -axis.

	Study	Treated Events	Treated Total N	Control Events	Control Total N	Risk ratio	Log risk ratio	Std Err	Variance	J	
1	Schmitz, 2007	7	78	13	76	0.525	-0.645	0.440	0.194		
2	Planer, 2011	23	75	25	76	0.932	-0.070	0.239	0.057		
3	Zellweger, 2005	117	501	36	166	1.077	0.074	0.168	0.028		
4	Eisenberg, 2013	49	183	43	194	1.208	0.189	0.182	0.033		
5	Tashkin, 2001	21	204	17	200	1.211	0.192	0.311	0.097		
6	SMK, 2001	26	143	20	143	1.300	0.262	0.273	0.074		
7	Witcher, 2011	22	108	27	175	1.320	0.278	0.260	0.068		
8	Nides, 2006	8	128	6	127	1.323	0.280	0.525	0.276		
9	Hurt, 1997	21	156	15	153	1.373	0.317	0.318	0.101		
10	Piper, 2007	42	224	21	156	1.393	0.331	0.246	0.061		
11	Jorenby, 2006	50	342	35	341	1.424	0.354	0.207	0.043		
12	Brown, 2007	38	255	27	269	1.485	0.395	0.236	0.056		
13	Rigotti, 2006	25	124	17	127	1.506	0.410	0.288	0.083		
14	Selby, 2003	18	141	12	143	1.521	0.420	0.353	0.125		
15	McCarthy, 2008	48	229	32	234	1.533	0.427	0.208	0.043		
16	Rovine, 2009	14	40	7	36	1.800	0.588	0.402	0.162		
17	Hall, 2002	13	73	7	73	1.857	0.619	0.439	0.192		
18	Fossati, 2007	101	400	26	193	1.874	0.628	0.202	0.041		
19	Tønnesen, 2003	111	527	20	180	1.896	0.640	0.227	0.052		
20	Gonzales, 2006	53	329	29	344	1.911	0.648	0.218	0.047		
21	Holt, 2005	19	88	5	46	1.986	0.686	0.463	0.220		
22	Ferry, 1994	13	95	6	95	2.167	0.773	0.472	0.223		
23	Tonstad, 2003	68	313	29	313	2.345	0.852	0.207	0.043		
24	Levine, 2010	42	195	12	156	2.800	1.030	0.309	0.096		
25	Jorenby, 1999	45	244	9	160	3.279	1.187	0.351	0.123		
26	Gonzales, 2001	20	226	5	224	3.965	1.377	0.491	0.241		
27	Ferry, 1992	10	23	0	22	20.125	3.002	1.418	2.010		
28											

Figure 50.13 Bupropion for smoking cessation. Data-entry screen.

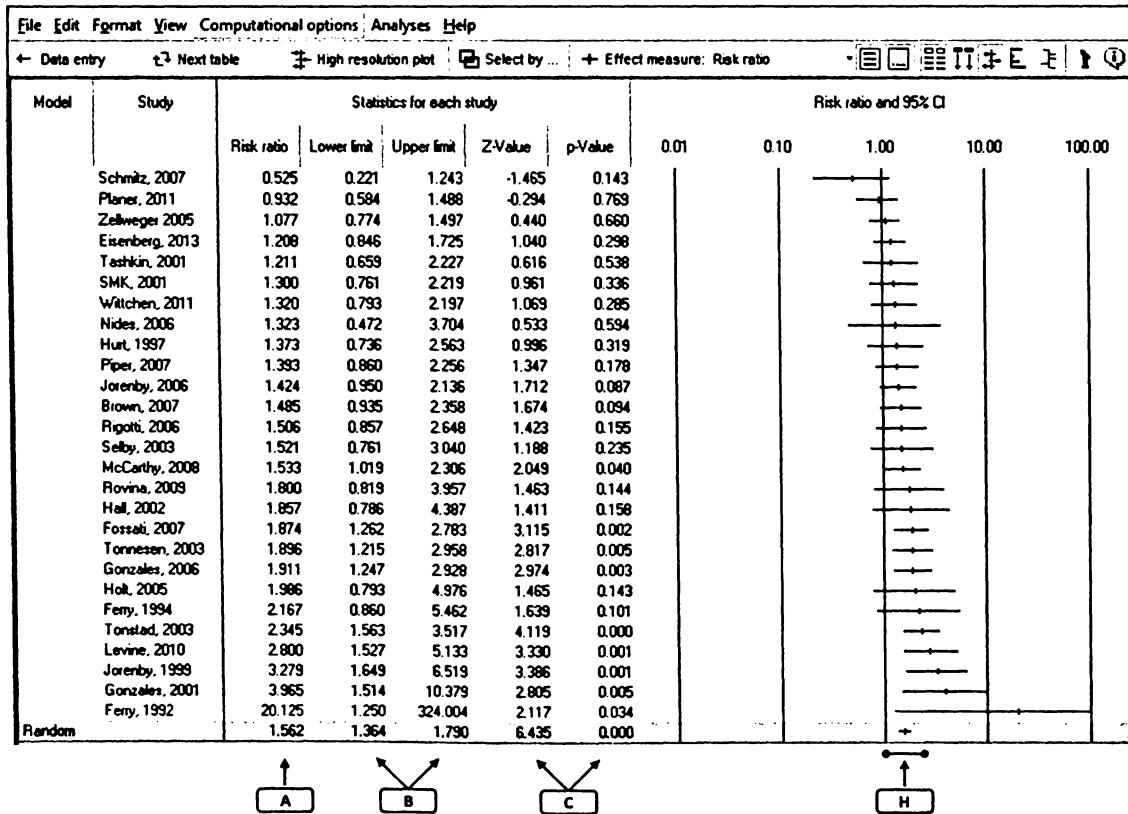
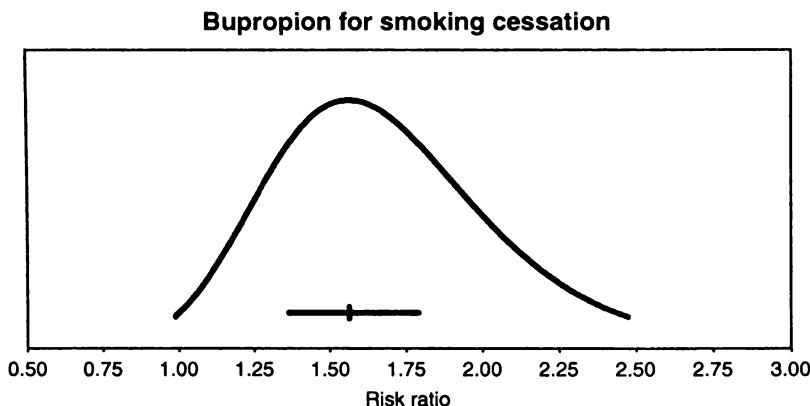


Figure 50.14 Bupropion for smoking cessation. Risk ratio &gt; 1 shows reduction in smoking.

Heterogeneity				Tau-squared			
Q-value	df [Q]	p-value	I-squared	Tau squared	Standard error	Variance	Tau
42.011	26	0.024	38.112	0.045	0.034	0.001	0.213

↓  
**D**  
 ↓  
**E**  
 ↓  
**F**  
 ↓  
**G**

Figure 50.15 Bupropion for smoking cessation. Heterogeneity statistics.



The mean effect size is 1.56 with a 95% confidence interval of 1.36 to 1.79  
 The true effect size in 95% of all comparable populations falls in the interval 0.99 to 2.47

Figure 50.16 Bupropion for smoking cessation. Distribution of true effects.

## MORTALITY FOLLOWING MITRAL-VALVE PROCEDURES IN ELDERLY PATIENTS

This example is a reanalysis of a systematic review published by Biancari *et al.* (2013). References in the text refer to Figure 50.17 through Figure 50.20.

### Overview

This analysis includes 25 studies of patients aged 80 years or more undergoing surgery for mitral valve repair or replacement. The outcome was the proportion of patients who died during or immediately following the surgery. The results of this analysis will be generalized to comparable studies, and so the random-effects model was employed for the analysis. The analysis was performed using a logit transformation of the data.

### What is the pooled estimate of mortality?

The pooled estimate of mortality is 0.166 [A], which means that *on average* about 17% of patients undergoing these procedures will not survive. The confidence interval for the prevalence is 0.142 to 0.193 [B], which tells us that the mean prevalence in the universe of comparable studies could fall anywhere in this range. Using the Knapp–Hartung adjustment, the 95% confidence interval is 0.134 to 0.203.

### How much does the effect size vary across studies?

The  $Q$ -statistic provides a test of the null hypothesis that all studies in the analysis share a common mortality risk. If all studies shared the same risk, the expected value of  $Q$  would be equal to the degrees of freedom (the number of studies minus 1). The  $Q$ -value is 65.668 with 24 degrees of freedom and  $p < 0.001$ . We can reject the null hypothesis that the true risk is identical in all these studies [D]. The  $I^2$  statistic is 63.453%, which tells us that some 64% of the variance in *observed* risk reflects variance in *true* risk rather than sampling error [E]. The variance of true risk ( $T^2$ ) is 0.080 in logit units [F], and the standard deviation of true risk ( $T$ ) is 0.282 in logit units [G]. If we assume that the prevalence is normally distributed (in logit units) we can estimate that the prediction interval is 0.097 to 0.269 [H] and predict that the true prevalence for any single population will usually fall in this range. The distribution of risk is shown in Figure 50.20.

The prediction interval is only a first step and tells us that the variation in mortality is of substantive import. What is clear from the prediction interval is that it would not be appropriate to focus on the mean prevalence of 18%, since this rate is not typical of many populations. Rather, we need to recognize that in any given population the prevalence could be substantially lower or higher than the mean. For example, the reviewers use meta-regression to see if the risk of death was lower in more recent studies (as techniques had improved) and if it was associated with the type of surgery (repair vs. replacement) or elements of surgery such as cross-clamp time.

Note that the distribution of effects is assumed to be symmetric in logit units. It appears to be skewed because the plot uses the risk rather than the logit on the  $x$ -axis.

	Study name	Events	Sample size	Event rate	Logit event rate	Std Err	Variance	H	I	J	
1	DiGregorio, 2004	1	59	0.017	-4.060	1.009	1.017				
2	Unic, 2005	6	106	0.057	-2.813	0.420	0.177				
3	Bhamidipati, 2011	6	81	0.074	-2.526	0.424	0.180				
4	Maleszka, 2008	5	55	0.091	-2.303	0.469	0.220				
5	Nlologe, 2012	12	129	0.093	-2.277	0.303	0.092				
6	Akins, 1997	4	42	0.095	-2.251	0.526	0.276				
7	Schmidler, 2008	2	21	0.095	-2.251	0.743	0.553				
8	Asimakopoulos,	10	96	0.104	-2.152	0.334	0.112				
9	Chicwe, 2011	43	322	0.134	-1.870	0.164	0.027				
10	Negndran, 2005	9	58	0.155	-1.695	0.363	0.132				
11	Ralph-Edwards,	3	18	0.167	-1.609	0.632	0.400				
12	Mehta, 2002	462	2720	0.170	-1.587	0.051	0.003				
13	Bessone, 2007	14	79	0.177	-1.535	0.295	0.087				
14	Craver, 1999	5	27	0.185	-1.482	0.495	0.245				
15	Aoyagi, 2010	4	21	0.190	-1.447	0.556	0.309				
16	Alexander, 2000	18	92	0.196	-1.414	0.263	0.069				
17	Jameieson, 1999	266	1351	0.197	-1.406	0.068	0.005				
18	Bashour, 1990	2	10	0.200	-1.386	0.791	0.625				
19	Tsai bis, 1994	7	31	0.226	-1.232	0.430	0.185				
20	Leavitt, 2009	35	147	0.238	-1.163	0.194	0.038				
21	Koh, 2001	3	12	0.250	-1.099	0.667	0.444				
22	Ngeage, 2008	4	16	0.250	-1.099	0.577	0.333				
23	Tsai, 1994	12	42	0.286	-0.916	0.342	0.117				
24	Freeman, 1991	10	27	0.370	-0.531	0.399	0.159				
25	Naunheim, 1990	5	10	0.500	0.000	0.632	0.400				
26											

Figure 50.17 Mortality following mitral-valve surgery in elderly patients. Data-entry screen.

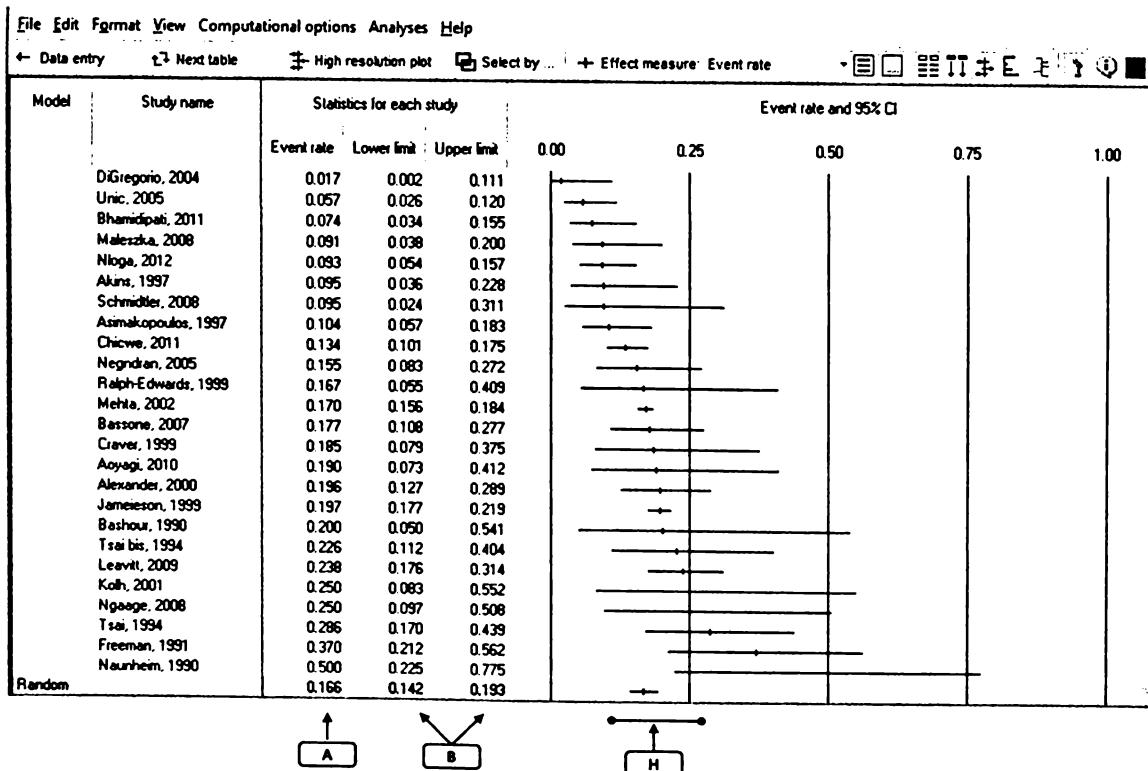
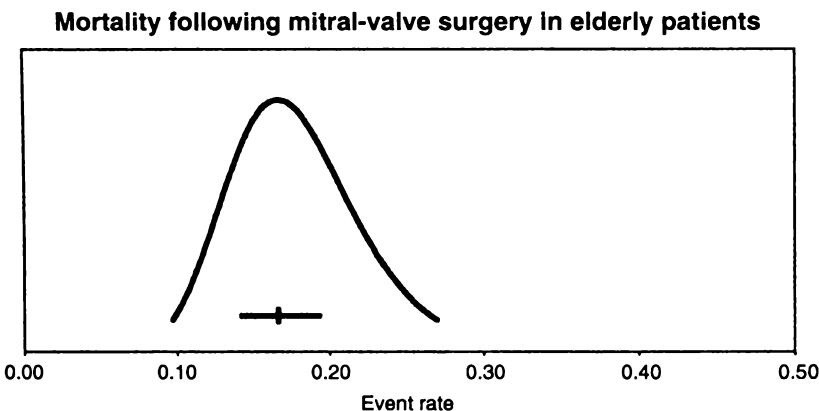


Figure 50.18 Mortality following mitral-valve surgery in elderly patients.

Heterogeneity				Tau-squared			
Q-value	df [Q]	p-value	I-squared	Tau squared	Standard error	Variance	Tau
65.668	24	0.000	63.453	0.080	0.070	0.005	0.282

↓  
 D      E      F      G  
 ↑  
 ↓  
 ↓  
 ↓

Figure 50.19 Mortality following mitral-valve surgery in elderly patients. Heterogeneity statistics.



The mean effect size is 0.17 with a 95% confidence interval of 0.14 to 0.19  
 The true effect size in 95% of all comparable populations falls in the interval 0.10 to 0.27

Figure 50.20 Mortality following mitral-valve surgery in elderly patients. Distribution of true risks.

# **Resources**



# Software for Meta-Analysis

---

Comprehensive meta-analysis

Metafor

Stata

Revman

---

### COMPREHENSIVE META-ANALYSIS

The screenshots in this volume are from the software Comprehensive Meta-Analysis (CMA), and Chapter 49 shows how to perform a meta-analysis using this software. CMA was initially released in 2000 and has been updated on a regular basis since then. The next version is scheduled for release in 2021.

For information on this program, visit [www.Meta-Analysis.com](http://www.Meta-Analysis.com).

The site allows you to download a free trial, and also features videos and worked examples.

Full disclosure. The authors of this volume are also the developers of CMA and some have a financial interest in the program.

### METAFOR

For researchers who would prefer to use R to perform meta-analysis, Wolfgang Viechtbauer has published a program called Metafor. To read more about the package, see

<http://www.metafor-project.org/>

<https://wviechtb.github.io/metafor/>

This is now complemented by code which allows the user to replicate the examples from the first edition of *Introduction to Meta-Analysis*. This may be downloaded at

[http://www.metafor-project.org/misc/borenstein2009.html#General\\_Notes\\_Setup](http://www.metafor-project.org/misc/borenstein2009.html#General_Notes_Setup)

## STATA

There are two options for running meta-analysis in Stata.

One option is to use macros that were published by various Stata users. The guide to these macros is *Meta-Analysis in Stata: An Updated Collection from the Stata Journal*, Second Edition, edited by Tom Palmer and Jonathan Sterne.

A second option is to use commands that are included in the most recent version of Stata. Information on these new commands is available at

<https://www.stata.com/features/meta-analysis/>

## REVMAN

Revman is the software employed by the Cochrane Collaboration.

For information on the latest release, see

<https://training.cochrane.org/online-learning/core-software-cochrane-reviews/revman>

# Web Sites, Societies, Journals, and Books

---

Web sites

Professional societies

Journals

Special issues dedicated to meta-analysis

Books on systematic review methods and meta-analysis

---

### WEB SITES

#### Introduction to meta-analysis

[www.Introduction-to-Meta-Analysis.com](http://www.Introduction-to-Meta-Analysis.com)

This is the website for this book.

- Download the data-sets used in this volume.
- Watch videos that illustrate some of the concepts in this book.
- Download software.

#### Comprehensive Meta-Analysis Software

[www.Meta-Analysis.com](http://www.Meta-Analysis.com)

This is the site for this book and for the program Comprehensive Meta-Analysis (CMA). The site includes the following:

- Free trial of Comprehensive Meta-Analysis.
- Worked examples for this book and for other books on meta-analysis.
- Papers on meta-analysis.
- Links to web sites on meta-analysis and systematic reviews.
- Links to courses and workshops on meta-analysis and systematic reviews.
- Links to relevant conferences.
- Full disclosure. Some of this book's authors have a financial interest in this software.

## Meta-Analysis Workshops

[www.Meta-Analysis-Workshops.com](http://www.Meta-Analysis-Workshops.com)

- A link to in-person and online workshops presented by Michael Borenstein.
- Full disclosure. Some of this book's authors have a financial interest in these workshops.

## The Human Genome Epidemiology Network

[www.cdc.gov/genomics/hugenet](http://www.cdc.gov/genomics/hugenet)

The Human Genome Epidemiology Network (HuGENet) is a global collaboration committed to the assessment of the impact of human genome variation on population health. Their website includes numerous systematic reviews (HuGE reviews) and other resources to support meta-analyses of genetic association studies.

## Hunter–Schmidt Meta-Analysis

<http://hunter-schmidt-meta-analysis.com>

This site offers information about the meta-analysis program developed by Hunter and Schmidt. It follows the Hunter–Schmidt approach to meta-analysis, with emphasis on artifact corrections.

## DARE, NHS EED, and HTA

<https://www.crd.york.ac.uk/CRDWeb/>

High-quality evidence to inform decision-making can be difficult to access, identify, and appraise. Our databases provide access to many thousands of systematic reviews, economic evaluations, and health technology assessments.

## Centre for Reviews and Dissemination

<https://www.york.ac.uk/crd/>

There is an ever-growing evidence base relating to the effectiveness and cost-effectiveness of healthcare interventions, but for clinicians and decision-makers this literature can be difficult and time-consuming to identify and appraise. Funded by NIHR, the CRD databases are providing the solution.

## Adam Hafdahl's Bibliography of Methods Papers

<https://adamhafdahl.net/bibliography/>

This is a free bibliography on methodology for research synthesis.

## The James Lind Library

[www.jameslindlibrary.org](http://www.jameslindlibrary.org)

This website is dedicated to 'Explaining and illustrating the development of fair tests of treatments in health care'. As such, it provides a history of the field, from

the earliest randomized trials to the most recent developments. Iain Chalmers and his colleagues have gathered key historical documents that help to put all current work in fascinating perspective.

## AHRQ

<https://srdrplus.ahrq.gov>

Software to facilitate data extraction.

## Brown School of Public Health

<https://www.brown.edu/public-health/cesh/resources/software>

Links to various open-source software.

## Robot Reviewer

<https://www.robotreviewer.net/about>

Software for automation of evidence synthesis.

## Using R

Wolfgang Viechtbauer

<http://www.metafor-project.org/>

<https://wviechtb.github.io/metafor/>

[https://github.com/wviechtb/meta\\_analysis\\_books](https://github.com/wviechtb/meta_analysis_books)

Mike Cheung (Structural Equation Modeling)

<https://cran.r-project.org/package=metaSEM>

Schwarzer, Carpenter, & Rucker

<https://meta-analysis-with-r.org>

## Meta-Analysis in Economics and Related Fields

<https://www.maer-net.org/>

[https://www.maer-net.org/post/towards-a-credibility-revolution.](https://www.maer-net.org/post/towards-a-credibility-revolution)

<https://www.deakin.edu.au/business/research/delmar>

## Meta-Analysis Based on Correlation Matrices in Social Science

<https://hubmeta.com/>

To extract data from correlation-matrix focused studies.

## Meta-Analysis for Ecology and Related Fields

<http://environmentalcomputing.net/meta-analysis/>

## PROFESSIONAL SOCIETIES

### The Cochrane Collaboration

<https://www.cochrane.org/>

The Cochrane Collaboration is an international not-for-profit organization, providing up-to-date information about the effects of health care. The organization's principal output is The Cochrane Library. This includes the Cochrane Database of Systematic Reviews, a regularly updated database of thousands of Cochrane reviews on the effects of healthcare interventions, on the accuracy of diagnostic tests, and on the methodology of systematic reviews and meta-analysis.

### The Campbell Collaboration

[www.campbellcollaboration.org](http://www.campbellcollaboration.org)

The international Campbell Collaboration (C2) is a nonprofit organization that aims to help people make well-informed decisions about the effects of interventions in the social, behavioral, and educational arenas.

### Society for Research Synthesis Methodology

<http://www.srsm.org/>

SRSM is a cross-disciplinary society that supports and promotes the development and use of innovative and robust methods of research synthesis.

## JOURNALS

### *Research Synthesis Methods*

<https://onlinelibrary.wiley.com/journal/17592887>

*Research Synthesis Methods* is a multidisciplinary peer-reviewed journal devoted to the development and dissemination of methods for designing, conducting, analyzing, interpreting, reporting, and applying systematic research synthesis.

### *Systematic Reviews*

<https://systematicreviewsjournal.biomedcentral.com/>

*Systematic Reviews* encompasses all aspects of the design, conduct, and reporting of systematic reviews. The journal aims to publish high-quality systematic review products including systematic review protocols, systematic reviews related to a very broad definition of health, rapid reviews, updates of already completed systematic reviews, and methods research related to the science of systematic reviews, such as decision modeling. The journal also aims to ensure that the results of all well-conducted systematic reviews are published, regardless of their outcome.

## SPECIAL ISSUES DEDICATED TO META-ANALYSIS

*Statistics in Medicine* 1987, vol. 6, no. 3: themed issue on meta-analysis.

*Journal of Educational and Behavioral Statistics* 1992, vol. 17, no. 4: special issue on meta-analysis.

*Statistical Methods in Medical Research* 1993, vol. 2, no. 2: themed issue on meta-analysis.

*International Journal of Epidemiology* 2002, vol. 11, no. 1: themed issue on systematic reviews and meta-analysis.

*Statistics in Medicine* 2002 vol. 21, no 11: proceedings from the 3rd Symposium on Systematic Review Methodology.

*Statistical Methods in Medical Research* 2001, vol. 10, no. 4: themed issue on meta-analysis, overviews and publication bias.

*Statistics in Medicine* 2002, volume 21, issue 11: Third Symposium on Systematic Review Methodology.

## Books on systematic review methods and meta-analysis

- Borenstein, M. (2019). *Common Mistakes in Meta-Analysis and How to Avoid Them*. Englewood, NJ: Biostat, Inc.
- Card, N.A. (2012). *Applied Meta-Analysis for Social Science Research*. New York, NY: The Guilford Press.
- Cheung, M.W.-L. (2015). *Meta-Analysis: A Structural Equation Modeling Approach*. Chichester, UK: Wiley.
- Cooper, H. (1998). *Synthesizing Research: A Guide for Literature Reviews* (3rd edn). Thousand Oaks, CA: Sage.
- Cooper, H. (2016). *Research Synthesis and Meta-Analysis: A Step-by-Step Approach* (5th edn). Thousand Oaks, CA: Sage.
- Cooper, H., Hedges, L.V., & Valentine, J.C. (2019). *The Handbook of Research Synthesis and Meta-Analysis* (3rd edn). New York: Russell Sage Foundation.
- Dias, S., Ades, A.E., Welton, N.J., Jansen, J.P., & Sutton, A.J. (2018). *Network Meta-analysis for Comparative Effectiveness Research*. Hoboken, NJ: Wiley.
- Egger, M., Higgins, J.P.T., Davey Smith, G. (forthcoming). *Systematic Reviews in Health Care: Meta-analysis in Context* (3rd edn). Chichester, UK: John Wiley & Sons.
- Glasziou, P., Irwig, L., Bain, C., & Colditz, G. (2001). *Systematic Reviews in Health Care: A Practical Guide*. Cambridge, UK: Cambridge University Press.
- Hartung, J., Knapp, G., & Sinha, B.K. (2008). *Statistical Meta-Analysis with Applications*. Hoboken, NJ: John Wiley & Sons.
- Hedges, L.V., & Olkin, I. (1985). *Statistical Methods for Meta-Analysis*. Orlando, FL: Academic Press.
- Higgins, J.P.T., Thomas, J., Chandler, J., Cumpston, M., Li, T., Page, M.J., Welch, V.A. (editors) (2019). *Cochrane Handbook for Systematic Reviews of Interventions*. 2nd Edition. Chichester, UK: John Wiley & Sons, Ltd.
- Huedo-Medina, T.B., & Johnson, B.T. (2010). *Modelos Estadísticos en Meta-análisis [Statistical Models in Meta-analysis]*. A Coruña, Spain: Netbiblio.
- Hunter, J.E., & Schmidt, F.L. (2015). *Methods of Meta-Analysis: Correcting Error and Bias in Research Findings* (3rd edn). Thousand Oaks, California: SAGE.
- Jak, S. (2015). *Meta-Analytic Structural Equation Modeling*. Switzerland: Springer International Publishing.

- Khan, K., Kunz, R., Kleijnen, J., & Antes, G. (2011). *Systematic Reviews to Support Evidence-Based Medicine*. Boca Raton, FL: CRC Press.
- Lipsey, M.W., & Wilson, D.B. (2001). *Practical Meta-Analysis*. Thousand Oaks, CA: Sage Publications.
- Littell, J.H., Corcoran, J., & Pillai, V. (2008). *Systematic Reviews and Meta-Analysis*. Oxford, UK: Oxford University Press.
- Neumann, P.J., Sanders, G.D., Russell, L.B., Siegel, J.E., & Ganiats, T.G. (eds) (2017). *Cost-Effectiveness in Health and Medicine* (2nd edn). New York, NY: Oxford University Press.
- Petticrew, M., & Roberts, H. (2006). *Systematic Reviews in the Social Sciences: A Practical Guide*. Oxford, UK: Blackwell.
- Pigott, T. (2012). *Advances in Meta-Analysis*. New York, NY: Springer.
- Rothstein, H., Sutton, A.J., & Borenstein, M. (2005). *Publication Bias in Meta-Analysis: Prevention, Assessment and Adjustments*. Chichester, England; Hoboken, NJ: Wiley.
- Schmid, C.H., Stijnen, T., & White, I.R. (eds) (2020). *Handbook of Meta-Analysis*. New York, NY: CRC Press.
- Schwarzer, G., Carpenter, J.R., & Rucker, G. (2015). *Meta-Analysis with R (Use R!)*. Switzerland: Springer International Publishing.
- Stangl, D.K., & Berry, D.A. (2000). *Meta-Analysis in Medicine and Health Policy*. New York, NY: Marcel Dekker.
- Sutton, A.J., Abrams, K.R., Jones, D.T., & Song, F. (2000). *Methods for Meta-Analysis in Medical Research*. Chichester, UK: John Wiley & Sons, Ltd.
- Welton, N.J., Sutton, A.J., Cooper, N.J., Abrams, K.R., & Ades, A.E. (2012). *Evidence Synthesis for Decision Making in Healthcare*. Chichester, UK: John Wiley and Sons.
- Whitehead, A. (2002). *Meta-Analysis of Controlled Clinical Trials*. Chichester, UK: John Wiley & Sons.

# References

---

- Abelson, R.P. (1997). A retrospective on the significance test ban of 1999 (If there were no significance tests, they would be invented). In: *What if There Were No Significance Tests?* (eds. L.L. Harlow, S.A. Mulaik and J.H. Steiger). Mahwah, NJ: Lawrence Erlbaum Associates.
- Arthur, W. Jr., Bennett, W. Jr., and Huffcutt, A.I. (2001). *Conducting Meta-Analysis Using SAS*. Mahwah, NJ: Lawrence Erlbaum Associates.
- Bailar, J.C. 3rd. (1997). The promise and problems of meta-analysis. *New England Journal of Medicine* 337: 559–561.
- Bax, L. and Moons, K.G. (2011). Beyond publication bias. *Journal of Clinical Epidemiology* 64 (5): 459–462. <https://doi.org/10.1016/j.jclinepi.2010.09.003>.
- Bayarri, M.J. (1988). Selection models and the file drawer problem: Comment. *Statistical Science* 3 (1): 128–131.
- Becker, B.J. (2005). Failsafe N or file-drawer number. In: *Publication Bias in Meta-Analysis: Prevention, Assessment and Adjustments* (eds. H.R. Rothstein, A.J. Sutton and M. Borenstein). Chichester, UK: John Wiley & Sons, Ltd.
- Begg, C.B. and Mazumdar, M. (1994). Operating characteristics of a rank correlation test for publication bias. *Biometrics* 50: 1088–1101.
- Berkey, C.S., Hoaglin, D.C., Mosteller, F., and Colditz, G.A. (1995). A random-effects regression model for meta-analysis. *Statistics in Medicine* 14: 395–411.
- Biancari, F., Schifano, P., Pighi, M. et al. (2013). Pooled estimates of immediate and late outcome of mitral valve surgery in octogenarians: a meta-analysis and meta-regression. *Journal of Cardio-thoracic and Vascular Anesthesia* 27 (2): 213–219. <https://doi.org/10.1053/j.jvca.2012.11.007>.
- Birge, R.T. (1941). The general physical constants. *Reports on Progress in Physics* 8: 90–101.
- Boccia, S., Hung, R., Ricciardi, G. et al. (2008). Meta- and pooled analyses of the methylenetetrahydrofolate reductase C677T and A1298C polymorphisms and gastric cancer risk: a huge-GSEC review. *American Journal of Epidemiology* 167: 505–516.
- Borenstein, M. (1994). The case for confidence intervals in controlled clinical trials. *Controlled Clinical Trials* 15: 411–428.
- Borenstein, M. (2000). The shift from significance testing to effect size estimation. In: *Comprehensive Clinical Psychology*, vol. 3 (eds. A.S. Bellack and M. Hersen), 313–349. Oxford, UK: Pergamon.
- Borenstein, M. (2019). *Common Mistakes in Meta-Analysis and How to Avoid Them*. Englewood, NJ: Biostat, Inc.
- Borenstein, M. (2020). Research Note: In a meta-analysis, the  $I^2$  index does not tell us how much the effect size varies across studies. *Journal of Physiotherapy* 66 (2): 135–139. <https://doi.org/10.1016/j.jphys.2020.02.011>.
- Borenstein, M., Hedges, L.V., Higgins, J.P.T., and Rothstein, H.R. (In preparation). *Meta-Regression – Multiple Regression in Meta-Analysis*. Englewood, NJ: Biostat, Inc.
- Borenstein, M., Hedges, L., Higgins, J.P.T., and Rothstein, H.R. (in preparation). *Computing Effect Sizes for Meta-analysis*. Chichester: John Wiley & Sons, Ltd.

- Borenstein, M., Higgins, J.P.T., Hedges, L.V., and Rothstein, H.R. (2017). Basics of meta-analysis:  $I^2$  is not an absolute measure of heterogeneity. *Research Synthesis Methods* 8 (1): 5–18. <https://doi.org/10.1002/jrsm.1230>.
- Butler, T.L. (1988). *The relationship of passive smoking to various health outcomes among Seventh-Day Adventists in California*. Los Angeles (Dissertation): University of California.
- Cannon, C.P., Steinberg, B.A., Murphy, S.A. et al. (2006). Meta-analysis of cardiovascular outcomes trials comparing intensive versus moderate statin therapy. *Journal of the American College of Cardiology* 48: 438–445.
- Card, N.A. (2012). *Applied Meta-Analysis for Social Science Research*. New York, NY: The Guilford Press.
- Carter, E.C., Schönbrodt, F.D., Gervais, W.M., and Hilgard, J. (2019). Correcting for bias in psychology: a comparison of meta-analytic methods. *Advances in Methods and Practices in Psychological Science* 2 (2): 115–144. <https://doi.org/10.1177/2515245919847196>.
- Castells, X., Ramos-Quiroga, J.A., Rigau, D. et al. (2011). Efficacy of methylphenidate for adults with attention-deficit hyperactivity disorder: a meta-regression analysis. *CNS Drugs* 25 (2): 157–169. <https://doi.org/10.2165/11539440-00000000-00000>.
- Chalmers, I. (2006). The scandalous failure of scientists to cumulate scientifically. Abstract to paper presented at: Ninth World Congress on Health Information and Libraries; 2005 Sep 20-23; Salvador, Brazil. (Available online: <http://www.icml9.org/program/activity.php?lang=pt&id=21>. Accessed on February 3, 2009).
- Chalmers, I. (2007). The lethal consequences of failing to make use of all relevant evidence about the effects of medical treatments: the need for systematic reviews. In P. Rothwell (ed.), *Treating Individuals: From Randomized Trials to Personalised Medicine* (pp. 37-58). London, UK: Elsevier.
- Chan, A.W., Hróbjartsson, A., Haahr, M.T. et al. (2004). Empirical evidence for selective reporting of outcomes in randomized trials: comparison of protocols to published articles. *Journal of the American Medical Association* 291: 2457–2465.
- Chan, W.C., and Fung, S.C. (1982). Lung cancer in non-smokers in Hong Kong. In: E. Grundmann (ed.), *Cancer Campaign*. (Vol 6. Cancer Epidemiology, pp. 199–202). New York, NY: Gustav Fischer.
- Chiolero, A., Santschi, V., Burnand, B. et al. (2012). Meta-analyses: with confidence or prediction intervals? *European Journal of Epidemiology* 27 (10): 823–825.
- Clarke, M. and Chalmers, I. (1998). Discussion sections in reports of controlled trials published in general medical journals: islands in search of continents? *Journal of the American Medical Association* 280: 280–282.
- Clarke, M. and Clarke, T. (2000). A study of the references used in Cochrane protocols and reviews. Three bibles, three dictionaries, and nearly 25,000 otherthings. *International Journal of Technology Assessment in Health Care* 16: 907–909.
- Cooper, H. (1998). *Synthesizing Research: A Guide for Literature Reviews*, 3rd. Thousand Oaks, CA: Sage.
- Cooper, H. (2016). *Research Synthesis and Meta-Analysis: A Step-by-Step Approach*, 5th. Thousand Oaks, CA: Sage.
- Cooper, H., Hedges, L.V., and Valentine, J.C. (2019). *The Handbook of Research Synthesis and Meta-Analysis*, 3rd. New York: Russell Sage Foundation.
- Coory, M.D. (2009). Comment on: heterogeneity in meta-analysis should be expected and appropriately quantified. *International Journal of Epidemiology* 39 (3): 932–932. <https://doi.org/10.1093/ije/dyp157>.
- Cohen, J. (1962). The statistical power of abnormal-social psychological research: a review. *Journal of Abnormal and Social Psychology* 65: 145–153.
- Cohen, J. (1969). *Statistical Power Analysis for the Behavioral Sciences*. New York, NY: Academic Press.
- Cohen, J. (1987). *Statistical Power Analysis for the Behavioral Sciences*. Hillside, NJ: Lawrence Erlbaum Associates.

- Colditz, G.A., Brewer, T.F., Berkey, C.S. et al. (1994). Efficacy of BCG vaccine in the prevention of tuberculosis. Meta-analysis of the published literature. *Journal of the American Medical Association* 271: 698–702.
- Crocker, L. and Algina, J. (1986). *Introduction to Classical and Modern Test Theory*. New York, NY: Holt, Rinehart, & Winston.
- Davey, J., Turner, R.M., Clarke, M.J., and Higgins, J.P.T. (2011). Characteristics of meta-analyses and their component studies in the Cochrane Database of Systematic Reviews: a cross-sectional, descriptive analysis. *BMC Med Res Methodol* 11: 160.
- Devereaux, P.J., Beattie, W.S., Choi, P.T. et al. (2005). How strong is the evidence for the use of perioperative beta blockers in non-cardiac surgery? Systematic review and meta-analysis of randomised controlled trials. *BMJ* 331: 313–321.
- Dickersin, K. (2005). Publication bias: Recognizing the problem, understanding its origins and scope, and preventing harm. In: *Publication Bias in Meta-Analysis* (eds. H.R. Rothstein, A. Sutton and M. Borenstein), 9–33. Chichester, UK: John Wiley & Sons, Ltd.
- Dickersin, K. and Min, Y.I. (1993a). NIH clinical trials and publication bias. *Online Journal of Current Clinical Trials*, Doc No 50.
- Dickersin, K., and Min, Y. I. (1993b). Publication bias: the problem that won't go away. *Annals of the New York Academy of Sciences*, 703, 135–146; discussion 146–148.
- Dickersin, K., Min, Y.I., and Meinert, C.L. (1992). Factors influencing publication of research results. Follow-up of applications submitted to two institutional review boards. *Journal of the American Medical Association* 267: 374–378.
- DuMouchel, W. (1988). Selection models and the file drawer problem: Comment. *Statistical Science* 3 (1): 132–133.
- Duval, S. and Tweedie, R. (2000a). A nonparametric 'trim and fill' method of accounting for publication bias in meta-analysis. *Journal of the American Statistical Association* 95: 89–98.
- Duval, S. and Tweedie, R. (2000b). Trim and fill: A simple funnel-plot-based method of testing and adjusting for publication bias in meta-analysis. *Biometrics* 56: 455–463.
- Easterbrook, P.J., Berlin, J.A., Gopalan, R., and Matthews, D.R. (1991). Publication bias in clinical research. *Lancet* 337: 867–872.
- Egger, M., Davey Smith, G., Schneider, M., and Minder, C. (1997). Bias in meta-analysis detected by a simple graphical test. *BMJ* 315: 629–634.
- Egger, M., Juni, P., Bartlett, C. et al. (2003). How important are comprehensive literature searches and the assessment of trial quality in systematic reviews? Empirical study. *Health Technol Assess* 7 (1): 1–76.
- Egger, M., Higgins, J.P.T., and Davey Smith, G. (forthcoming). *Systematic Reviews in Health Care: Meta-analysis in Context* (3rd edition). Chichester, UK: John Wiley & Sons.
- Eysenck, H.J. (1978). An exercise in mega-silliness. *American Psychologist* 33: 517–519.
- Feinstein, A.R. (1995). Meta-analysis: statistical alchemy for the 21st century. *Journal of Clinical Epidemiology* 48: 71–79.
- Freiman, J.A., Chalmers, T.C., Smith, Jr., H., and Kuebler, R.R. (1978). The importance of beta, the type II error and sample size in the design and interpretation of the randomized control trial. Survey of 71 'negative' trials. *New England Journal of Medicine*, 299, 690–694.
- Gilbert, R., Salanti, G., Harden, M., and See, S. (2005). Infant sleeping position and the sudden infant death syndrome: systematic review of observational studies and historical review of recommendations from 1940 to 2002. *International Journal of Epidemiology* 34: 874–887.
- Giustini, D. (2019). Retrieving grey literature, information, and data in the digital age. In: *The Handbook of Research Synthesis and Meta-Analysis*, 3rd (eds. H. Cooper, L.V. Hedges and J.C. Valentine), 101–126. New York, NY: Russell Sage Foundation.
- Glanville, J. (2019). Searching bibliographic databases. In: *The Handbook of Research Synthesis and Meta-Analysis*, 3rd (eds. H. Cooper, L.V. Hedges and J.C. Valentine), 73–100. New York, NY: Russell Sage Foundation.
- Glasziou, P., Irwig, L., Bain, C., and Colditz, G. (2001). *Systematic Reviews in Health Care: A Practical Guide*. Cambridge, UK: Cambridge University Press.

- Glasziou, P.P. and Irwig, L.M. (1995). An evidence based approach to individualising treatment. *BMJ (Clinical research ed.)* 311 (7016): 1356–1359.
- Goodman, S.N. (1991). Have you ever meta-analysed something you didn't like? *Annals Internal Medicine* 114: 244–246.
- Gøtzsche, P.C. (1987). Reference bias in reports of drug trials. *BMJ* 295: 654–656.
- Graham, P.L. and Moran, J.L. (2012). Robust meta-analytic conclusions mandate the provision of prediction intervals in meta-analysis summaries. *Journal of Clinical Epidemiology* 65 (5): 503–510. <https://doi.org/10.1016/j.jclinepi.2011.09.012>.
- Grissom, R.J. and Kim, J.J. (2005). *Effect Sizes for Research: A Broad Practical Approach*. Mahwah, NJ: Lawrence Erlbaum Associates.
- Guddat, C., Grouwen, U., Bender, R., and Skipka, G. (2012). A note on the graphical presentation of prediction intervals in random-effects meta-analyses. *Systematic Reviews* 1: 34. <https://doi.org/10.1186/2046-4053-1-34>.
- Hackshaw, A.K., Law, M.R., and Wald, N.J. (1997). The accumulated evidence on lung cancer and environmental tobacco smoke. *BMJ* 315: 980–988.
- Halpern, S.D. and Berlin, J.A. (2005). Beyond conventional publication bias: other determinants of data suppression. In: *Publication Bias in Metaanalysis: Prevention, Assessment and Adjustments* (eds. H.R. Rothstein, A.J. Sutton and M. Borenstein). Chichester, UK: John Wiley & Sons, Ltd.
- Hartung, J., Cottrell, J.E., and Giffin, J.P. (1983). Absence of evidence is not evidence of absence. *Anesthesiology* 58: 298–300.
- Hartung, J., Knapp, G., and Sinha, B.K. (2008). *Statistical Meta-Analysis with Applications*. Hoboken, NJ: John Wiley & Sons, Inc.
- Hasselblad, V. and Hedges, L.V. (1995). Meta-analysis of screening and diagnostic tests. *Psychological Bulletin* 117: 167–178.
- Hedges, L. (1981). Distribution theory for Glass's estimator of effect size and related estimators. *Journal of Educational Statistics* 6: 107–128.
- Hedges, L. (1984). Estimation of effect size under nonrandom sampling: The effects of censoring studies yielding statistically insignificant mean differences. *Journal of Educational Statistics* 9: 61–85.
- Hedges, L. (1989) Estimating the normal mean and variance under a selection model. In L. Gleser, M.D. Perlman, S.J. Press, A.R. Sampson. *Contributions to Probability and Statistics: Essays in Honor of Ingram Olkin* (pp. 447–458). New York, NY: Springer Verlag.
- Hedges, L. and Olkin, I. (1980). Vote-counting methods in research synthesis. *Psychological Bulletin*. 88: 359–369.
- Hedges, L. and Olkin, I. (1985). *Statistical Methods for Meta-analysis*. San Diego, CA: Academic Press.
- Hedges, L. and Pigott, T.D. (2001). The power of statistical tests in meta-analysis. *Psychological Methods* 6: 203–217.
- Hedges, L. and Pigott, T.D. (2004). The power of statistical tests for moderators in meta-analysis. *Psychological Methods* 9: 426–445.
- Hedges, L., Gurevitch, J., and Curtis, P. (1999). The meta-analysis of response ratios in experimental ecology. *Ecology* 80: 1150–1156.
- Hedges, L.V. (1988). Selection models and the file drawer problem: Comment. *Statistical Science* 3 (1): 118–120.
- Hedges, L.V. (1992). Modeling publication selection effects in meta-analysis. *Statistical Science* 7 (2): 246–255. <https://doi.org/10.1214/ss/1177011364>.
- Hedges, L.V. and Vevea, J.L. (1998). Fixed and random-effects models in meta-analysis. *Psychological Methods* 3 (4): 486–504.
- Hedges, L.V. and Schauer, J.M. (2019). More Than One Replication Study Is Needed for Unambiguous Tests of Replication. *Journal of Educational and Behavioral Statistics*, 44 (5): 543–570.
- Higgins, J.P.T. (2008). Commentary: Heterogeneity in meta-analysis should be expected and appropriately quantified. *International Journal of Epidemiology* 37 (5): 1158–1160. <https://doi.org/10.1093/ije/dyn204>.

- Higgins, J.P.T., Thompson, S.G., and Spiegelhalter, D.J. (2009). A re-evaluation of random-effects meta-analysis. *Journal of the Royal Statistical Society. Series A, (Statistics in Society)* 172 (1): 137–159. <https://doi.org/10.1111/j.1467-985X.2008.00552.x>.
- Higgins, J.P.T. and Thompson, S.G. (2002). Quantifying heterogeneity in a meta-analysis. *Statistics in Medicine* 21: 1539–1558.
- Higgins, J.P.T., Thompson, S.G., Deeks, J.J., and Altman, D.G. (2003). Measuring inconsistency in meta-analyses. *BMJ* 327: 557–560.
- Higgins, J.P.T. and Thompson, S.G. (2004). Controlling the risk of spurious findings from metaregression. *Statistics in Medicine* 23: 1663–1682.
- Higgins, J.P.T., Thomas, J., Chandler, J., Cumpston, M., Li, T., Page, M.J., Welch, V.A. (editors) (2019). *Cochrane Handbook for Systematic Reviews of Interventions*. 2nd Edition. Chichester, UK: John Wiley & Sons, Ltd.
- Hjetland, H.N., Brinchmann, E.I., Scherer, R., and Melby-Lervåg, M. (2017). Preschool predictors of later reading comprehension ability: a systematic review. *Campbell Systematic Reviews* 14: 1–156. <https://doi.org/10.4073/csr.2017.14>.
- Holst, L.B., Petersen, M.W., Haase, N. et al. (2015). Restrictive versus liberal transfusion strategy for red blood cell transfusion: systematic review of randomised trials with meta-analysis and trial sequential analysis. *BMJ* 350: h1354. <https://doi.org/10.1136/bmj.h1354>.
- Hopewell, S., Clarke, M., and Mallett, S. (2005). Grey literature and systematic reviews. In: *Publication Bias in Meta-Analysis: Prevention, Assessment and Adjustments* (eds. H.R. Rothstein, A.J. Sutton and M. Bornstein). Chichester, UK: John Wiley & Sons, Ltd.
- Huedo-Medina, T.B. and Johnson, B.T. (2010). *Modelos Estadísticos en Meta-análisis [Statistical Models in Meta-analysis]*. In: *A Coruña*. Spain: Netbiblio.
- Huedo-Medina, T.B., Sanchez-Meca, J., Marin-Martinez, F., and Botella, J. (2006). Assessing heterogeneity in meta-analysis: Q statistic or I<sup>2</sup> index? *Psychol Methods* 11 (2): 193–206. <https://doi.org/10.1037/1082-989X.11.2.193>.
- Hughes, J.R., Stead, L.F., Hartmann-Boyce, J. et al. (2014). Antidepressants for smoking cessation. *Cochrane Database of Systematic Reviews* 1 <https://doi.org/10.1002/14651858.CD000031.pub4>.
- Hunter, J.E. and Schmidt, F.L. (1990). *Methods of Meta-analysis: Correcting Error and Bias in Research Findings*. Newbury Park, CA: Sage Publications.
- Hunter, J.E. and Schmidt, F.L. (2004). *Methods of Meta-analysis: Correcting Error and Bias in Research Findings*, 2nde. Newbury Park, CA: Sage Publications.
- Hunter, J.E. and Schmidt, F.L. (2015). *Methods of Meta-Analysis : Correcting Error and Bias in Research Findings*, 3rde. Thousand Oaks, CA: SAGE.
- IntHout, J., Ioannidis, J.P., and Borm, G.F. (2014). The Hartung-Knapp-Sidik-Jonkman method for random effects meta-analysis is straightforward and considerably outperforms the standard DerSimonian-Laird method. *BMC Medical Research Methodology* 14 (1): 25. <https://doi.org/10.1186/1471-2288-14-25>.
- IntHout, J., Ioannidis, J.P.A., Rovers, M.M., and Goeman, J.J. (2016). Plea for routinely presenting prediction intervals in meta-analysis. *BMJ Open* 6 (7): e010247. <https://doi.org/10.1136/bmjopen-2015-010247>.
- Ioannidis, J.P. (2008a). Interpretation of tests of heterogeneity and bias in meta-analysis. *Journal of Evaluation in Clinical Practice* 14 (5): 951–957. <https://doi.org/10.1111/j.1365-2753.2008.00986.x>.
- Ioannidis, J.P. (2008b). Why most discovered true associations are inflated. *Epidemiology* 19 (5): 640–648. <https://doi.org/10.1097/EDE.0b013e31818131e7>.
- Ioannidis, J.P. and Trikalinos, T.A. (2007). The appropriateness of asymmetry tests for publication bias in meta-analyses: a large survey. *Canadian Medical Association Journal* 176 (8): 1091–1096. <https://doi.org/10.1503/cmaj.060410>.
- Iyengar, S. and Greenhouse, J.B. (1988a). Selection Models and the File Drawer Problem. *Statistical Science* 3 (1): 109–117.
- Iyengar, S. and Greenhouse, J.B. (1988b). Selection models and the file drawer problem: Rejoinder. *Statistical Science* 3 (1): 133–135.

- Jackson, D., Law, M., Rücker, G., and Schwarzer, G. (2017). The Hartung-Knapp modification for random-effects meta-analysis: a useful refinement but are there any residual concerns? *Statistics in Medicine* 36 (25): 3923–3934. <https://doi.org/10.1002/sim.7411>.
- Juhl, C., Christensen, R., Roos, E.M. et al. (2014). Impact of exercise type and dose on pain and disability in knee osteoarthritis: a systematic review and meta-regression analysis of randomized controlled trials. *Arthritis & Rheumatology* 66 (3): 622–636. <https://doi.org/10.1002/art.38290>.
- Jüni, P., Holenstein, F., Sterne, J. et al. (2002). Direction and impact of language bias in meta-analyses of controlled trials: empirical study. *International Journal of Epidemiology* 31: 115–123.
- Kane, J.M. and Borenstein, M. (1985). Compliance in the long-term treatment of schizophrenia. *Psychopharmacology Bulletin* 21: 23–27.
- Katout, M., Zhu, H., Rutsky, J. et al. (2014). Effect of GLP-1 mimetics on blood pressure and relationship to weight loss and glycemia lowering: results of a systematic meta-analysis and meta-regression. *American Journal of Hypertension* 27 (1): 130–139. <https://doi.org/10.1093/ajh/hpt196>.
- Keith, B.G.D. and Begg, C.B. (1992). An approach for assessing publication bias prior to performing a meta-analysis. *Statistical Science* 7 (2): 237–245.
- Khan, K., Kunz, R., Kleijnen, J., and Antes, G. (2011). *Systematic Reviews to Support Evidence-Based Medicine*. Boca Raton, FL: CRC Press.
- Knapp, G. and Hartung, J. (2003). Improved tests for a random effects meta-regression with a single covariate. *Statistics in Medicine* 22: 2693–2710.
- Lacasaña-Navarro, M., Galvan-Portillo, M., Chen, J. et al. (2006). Methylenetetrahydrofolate reductase 677C > T polymorphism and gastric cancer susceptibility in Mexico. *European Journal of Cancer* 42: 528–533.
- Laird, N., Patil, G.P., and Taillie, C. (1988). Selection models and the file drawer problem: Comment. *Statistical Science* 3 (1): 126–128.
- Lau, J., Antman, E.M., Jimenez-Silva, J. et al. (1992). Cumulative meta-analysis of therapeutic trials for myocardial infarction. *New England Journal of Medicine* 327: 248–254.
- Lau, J. and Chalmers, T.C. (1995). The rational use of therapeutic drugs in the 21st century. Important lessons from cumulative meta-analyses of randomized control trials. *International Journal of Technology Assessment in Health Care* 11: 509–522.
- Lau, J., Schmid, C.H., and Chalmers, T.C. (1995). Cumulative meta-analysis of clinical trials builds evidence for exemplary medical care. *Journal of Clinical Epidemiology* 48: 45–57. discussion 59–60.
- Lefebvre, C., Glanville, J., Briscoe, S. et al. (2019). Searching for and selecting studies. In J.P.T. Higgins, J. Thomas, M. Cumpston, T. Li, M.J. Page, & V.A. Welch (eds), *Cochrane Handbook for Systematic Reviews of Interventions*, (2nd edition, pp. 67–108). Chichester, UK: John Wiley and Sons, Ltd.
- LeLorier, J., Gregoire, G., Benhaddad, A. et al. (1997). Discrepancies between meta-analyses and subsequent large randomized controlled trials. *New England Journal of Medicine* 337: 536–542.
- Lewis, S. and Clarke, M. (2001). Forest plots: trying to see the wood and the trees. *BMJ* 322 (7300): 1479–1480.
- Light, R.J. and Pillemer, D.B. (1984). *Summing up: The Science of Reviewing Research*. Cambridge, MA: Harvard University Press.
- Light, R.J., Singer, J.D., and Willett, J.B. (1994). The visual presentation and interpretation of meta-analyses. In: *The Handbook of Research Synthesis* (eds. M. Cooper and L.V. Hedges). New York, NY: Russell Sage Foundation.
- Linde, K., Clausius, N., Ramirez, G. et al. (1997). Are the clinical effects of homeopathy placebo effects? A meta-analysis of placebo-controlled trials. *Lancet* 350 (9081): 834–843.
- Linde, K., Scholz, M., Ramirez, G. et al. (1999). Impact of study quality on outcome in placebo-controlled trials of homeopathy. *Journal of Clinical Epidemiology* 52 (7): 631–636.
- Lipsey, M.W. and Wilson, D.B. (2001). *Practical Meta-Analysis*. Thousand Oaks, CA: Sage Publications.

- Littell, J.H., Corcoran, J., and Pillai, V. (2008). *Systematic Reviews and Meta-Analysis*. Oxford, UK: Oxford University Press.
- Lord, F.M. and Novick, M.R. (1968). *Statistical Theories of Mental Test Scores*. Reading, MA: Addison-Wesley.
- Mallet, S., Hopewell, S., and Clarke, M. (2002). The use of grey literature in the first 1000 Cochrane reviews. Paper presented at the Fourth Symposium on Systematic Reviews: Pushing the Boundaries; 2002 Jul, 2–4. Oxford, UK.
- Mann, C. (1990). Meta-analysis in the breach. *Science* 249: 476–480.
- McShane, B.B., Bockenholt, U., and Hansen, K.T. (2016). Adjusting for publication bias in meta-analysis: an evaluation of selection methods and some cautionary notes. *Perspectives on Psychological Science* 11 (5): 730–749. <https://doi.org/10.1177/1745691616662243>.
- Meehl, P.E. (1978). Theoretical risks and tabular asterisks: Sir Karl, Sir Ronald, and the slow progress of soft psychology. *Journal of Consulting and Clinical Psychology* 46: 806–834.
- Meehl, P.E. (1990). Why summaries of research on psychological theories are often uninterpretable. *Psychological Reports* 66: 195–244.
- Mittlböck, M. and Heinzl, H. (2006). A simulation study comparing properties of heterogeneity measures in meta-analyses. *Statistics in Medicine* 25 (24): 4321–4333. <https://doi.org/10.1002/sim.2692>.
- Moreno, S.G., Sutton, A.J., Ades, A.E. et al. (2009). Assessment of regression-based methods to adjust for publication bias through a comprehensive simulation study. *BMC Medical Research Methodology* 9: 2. <https://doi.org/10.1186/1471-2288-9-2>.
- Moreno, S.G., Sutton, A.J., Ades, A.E. et al. (2011). Adjusting for publication biases across similar interventions performed well when compared with gold standard data. *Journal of Clinical Epidemiology* 64 (11): 1230–1241. <https://doi.org/10.1016/j.jclinepi.2011.01.009>.
- Moreno, S.G., Sutton, A.J., Thompson, J.R. et al. (2012). A generalized weighting regression-derived meta-analysis estimator robust to small-study effects and heterogeneity. *Statistics in Medicine* 31 (14): 1407–1417. <https://doi.org/10.1002/sim.4488>.
- Moreno, S.G., Sutton, A.J., Turner, E.H. et al. (2009). Novel methods to deal with publication biases: secondary analysis of antidepressant trials in the FDA trial registry database and related journal publications. *BMJ* b2981: 339. <https://doi.org/10.1136/bmj.b2981>.
- Nagashima, K., Noma, H., and Furukawa, T.A. (2019). Prediction intervals for random-effects meta-analysis: A confidence distribution approach. *Statistical Methods in Medical Research* 28 (6): 1689–1702. <https://doi.org/10.1177/0962280218773520>.
- Normand, S.L. (1999). Meta-analysis: formulating, evaluating, combining, and reporting. *Statistics in Medicine* 18: 321–359.
- O'Farrell, N. and Egger, M. (2000). Circumcision in men and the prevention of HIV infection: a 'meta-analysis' revisited. *International Journal of STD and AIDS* 11: 137–142.
- Orwin, R.G. and Boruch, R.F. (1983). RRTmeets RDD: statistical strategies for assuring response privacy in telephone surveys. *Public Opinion Quarterly* 46: 560–571.
- Page, M.J., Higgins, J.P.T., and Sterne, J.A.C. (2019). Assessing risk of bias due to missing results in a synthesis. In J.P.T. Higgins, J. Thomas, M. Cumpston, T. Li, M.J. Page, & V.A. Welch (eds), *Cochrane Handbook for Systematic Reviews of Interventions* (2nd edition, pp. 349–374). Chichester, UK: John Wiley and Sons, Ltd.
- Peters, J.L., Sutton, A.J., Jones, D.R. et al. (2007). Performance of the trim and fill method in the presence of publication bias and between-study heterogeneity. *Statistics in Medicine* 26 (25): 4544–4562. <https://doi.org/10.1002/sim.2889>.
- Peters, J.L., Sutton, A.J., Jones, D.R. et al. (2008). Contour-enhanced meta-analysis funnel plots help distinguish publication bias from other causes of asymmetry. *Journal of Clinical Epidemiology* 61 (10): 991–996. <https://doi.org/10.1016/j.jclinepi.2007.11.010>.
- Peters, J.L., Sutton, A.J., Jones, D.R. et al. (2010). Assessing publication bias in meta-analyses in the presence of between-study heterogeneity. *Journal of the Royal Statistical Society: Series A (Statistics in Society)* 173 (3): 575–591.

- Peto, R. (1987). Why do we need systematic overviews of randomized trials? (Transcript of an oral presentation, modified by the editors). *Statistics in Medicine* 6 (3): 233–240. <https://doi.org/10.1002/sim.4780060306>.
- Petticrew, M. and Roberts, H. (2006). *Systematic reviews in the social sciences: a practical guide*. Oxford, UK: Blackwell.
- Pigott, T. (2012). *Advances in Meta-Analysis*. New York, NY: Springer.
- Phillips, W.C., Scott, J.A., and Blasczynski, G. (1983). Statistics for diagnostic procedures. II. The significance of 'no significance': what a negative statistical test really means. *American Journal of Roentgenology* 141: 203–206.
- Pogue, J. and Yusuf, S. (1998). Overcoming the limitations of current meta-analysis of randomised controlled trials. *Lancet* 351: 47–52.
- Poole, C. and Greenland, S. (1999). Random-effects meta-analyses are not always conservative. *American Journal of Epidemiology* 150 (5): 469–475.
- Rao, C.R. (1988). Selection models and the file drawer problem: Comment. *Statistical Science* 3 (1): 131–131.
- Ravnskov, U. (1992). Frequency of citation and outcome of cholesterol lowering trials. *BMJ* 305: 717.
- Reed, J.F. 3rd and Slaichert, W. (1981). Statistical proof in inconclusive 'negative' trials. *Archives of Internal Medicine* 141: 1307–1310.
- Reed, J.G. and Baxter, P.M. (2009). Using reference databases. In: *The Handbook of Research Synthesis*, 2nde (eds. H. Cooper, L.V. Hedges and J. Valentine). New York, NY: Sage Publications.
- Reynolds, T.B. (1980). Type II error in clinical trials (editor's reply to letter). *Gastroenterology* 79: 180.
- Rice, K., Higgins, J.P.T., and Lumley, T. (2018). A re-evaluation of fixed effect(s) meta-analysis. *Journal of the Royal Statistical Society: Series A (Statistics in Society)* 181 (1): 205–227. <https://doi.org/10.1111/rssa.12275>.
- Riley, R.D., Higgins, J.P.T., and Deeks, J.J. (2011). Interpretation of random effects meta-analyses. *BMJ* 342: d549. <https://doi.org/10.1136/bmj.d549>.
- Rona, R.J., Keil, T., Summers, C. et al. (2007). The prevalence of food allergy: a meta-analysis. *Journal of Allergy and Clinical Immunology* 120: 638–646.
- Ronksley, P.E., Brien, S.E., Turner, B.J. et al. (2011). Association of alcohol consumption with selected cardiovascular disease outcomes: a systematic review and meta-analysis. *BMJ* 342: d671. <https://doi.org/10.1136/bmj.d671>.
- Rosenthal, R. (1979). The File drawer problem and tolerance for null results. *Psychological Bulletin* 86: 638–641.
- Rosenthal, R. and Rubin, D.B. (1988). Selection models and the file drawer problem: Comment: assumptions and procedures in the file drawer problem. *Statistical Science* 3 (1): 120–125.
- Rossi, J. (1997). A case study in the failure of psychology as a cumulative science: The spontaneous recovery of verbal learning. In: *What if There Were No Significance Tests?* (eds. L.L. Harlow, S.A. Mulaik and J.H. Steiger), 175–198. Mahwah, NJ: Laurence Erlbaum Associates.
- Rossi, J.S. (1990). Statistical power of psychological research: what have we gained in 20 years? *Journal of Consulting and Clinical Psychology* 58: 646–656.
- Rothstein, H.R. (2006). Use of unpublished data in systematic reviews in the Psychological Bulletin 1995–2005. Unpublished manuscript.
- Rothstein, H.R. and Hopewell, S. (2009). The Grey literature. In: *The Handbook of Research Synthesis*, 2nde (eds. H. Cooper, L.V. Hedges and J. Valentine). New York, NY: Sage Publications.
- Rothstein, H., Sutton, A.J., and Borenstein, M. (2005). *Publication Bias in Meta-Analysis: Prevention, Assessment and Adjustments*. Chichester. In: *England*. Hoboken, NJ: John Wiley & Sons, Inc.

- Rucker, G., Carpenter, J.R., and Schwarzer, G. (2011). Detecting and adjusting for small-study effects in meta-analysis. *Biometrical Journal* 53 (2): 351–368. <https://doi.org/10.1002/bimj.201000151>.
- Rucker, G., Schwarzer, G., Carpenter, J.R., and Schumacher, M. (2008). Undue reliance on I<sup>2</sup> in assessing heterogeneity may mislead. *BMC Medical Research Methodology* 8: 79. <https://doi.org/10.1186/1471-2288-8-79>.
- Rucker, G., Schwarzer, G., Carpenter, J.R. et al. (2011). Treatment-effect estimates adjusted for small-study effects via a limit meta-analysis. *Biostatistics* 12 (1): 122–142. <https://doi.org/10.1093/biostatistics/kxq046>.
- Sanchez-Meca, J., Marin-Martinez, F., and Chacon-Moscoso, S. (2003). Effect-size indices for dichotomized outcomes in meta-analysis. *Psychological Methods* 8: 448–467.
- Schmidt, F.L. (1996). Statistical significance testing and cumulative knowledge in psychology: Implications for training of researchers. *Psychological Methods* 1: 115–129.
- Sedlmeier, P. and Gigerenzer, G. (1989). Do studies of statistical power have an effect on the power of studies? *Psychological Bulletin* 105: 309–316.
- Shapiro, S. (1994). Meta-analysis/shmeta-analysis. *American Journal of Epidemiology* 140: 771–778.
- Sidik, K., and Jonkman, J.N. (2002). A simple confidence interval for meta-analysis. *Statistics in Medicine*, 21(21), 3153–3159. doi:[doi:10.1002/sim.1262](https://doi.org/10.1002/sim.1262)
- Simonsohn, U., Nelson, L.D., and Simmons, J.P. (2014). p-Curve and effect size: correcting for publication bias using only significant results. *Perspectives on Psychological Science* 9 (6): 666–681. <https://doi.org/10.1177/1745691614553988>.
- Sirmans, G.S., Macdonald, L., Macpherson, D.A., and Zietz, E.N. (2006). The value of housing characteristics: a meta analysis. *Journal of Real Estate Finance and Economics* 33: 215–240.
- Smith, G.D. and Egger, M. (1994). Who benefits from medical interventions? *BMJ* 308 (6921): 72–74.
- Sorita, A., Ahmed, A., Starr, S.R. et al. (2014). Off-hour presentation and outcomes in patients with acute myocardial infarction: systematic review and meta-analysis. *BMJ* f7393: 348. <https://doi.org/10.1136/bmj.f7393>.
- Stangl, D.K. and Berry, D.A. (2000). *Meta-Analysis in Medicine and Health Policy*. New York, NY: Marcel Dekker.
- Stanley, T.D. and Doucouliagos, H. (2014). Meta-regression approximations to reduce publication selection bias. *Research Synthesis Methods* 5 (1): 60–78. <https://doi.org/10.1002/jrsm.1095>.
- Sterne, J., Egger, M., and Moher, D. (eds.) (2008). *Addressing Reporting Biases*. Chichester, UK: John Wiley & Sons, Ltd.
- Sterne, J.A. and Egger, M. (2001). Funnel plots for detecting bias in meta-analysis: guidelines on choice of axis. *Journal of Clinical Epidemiology* 54: 1046–1055.
- Sterne, J.A., Egger, M., and Smith, G.D. (2001a). Systematic reviews in health care: Investigating and dealing with publication and other biases in meta-analysis. *BMJ* 323: 101–105.
- Sterne, J.A., Gavaghan, D., and Egger, M. (2000). Publication and related bias in meta-analysis: power of statistical tests and prevalence in the literature. *Journal of Clinical Epidemiology* 53: 1119–1129.
- Sterne, J.A., Sutton, A.J., Ioannidis, J.P. et al. (2011). Recommendations for examining and interpreting funnel plot asymmetry in meta-analyses of randomised controlled trials. *BMJ* d4002: 343. <https://doi.org/10.1136/bmj.d4002>.
- Sterne, J.A.C., Egger, M., and Davey Smith, G. (2001b). Investigating and dealing with publication and other biases. In: *Systematic Reviews in Health-Care: Meta-Analysis in Context* (eds. M. Egger, G. Davey-Smith and D. Altman), 189–208. London, UK: BMJ.

- Stewart, L.A. and Tierney, J.F. (2002). To IPD or not to IPD? Advantages and disadvantages of systematic reviews using individual patient data. *Evaluation and the Health Professions* 25: 76–97.
- Stuck, A.E., Rubenstein, L.Z., and Wieland, G.D. (1998). Asymmetry detected in funnel plot was probably due to true heterogeneity: bias in meta-analysis detected by a simple graphical test. *BMJ* 316 (7129): 469. <https://doi.org/10.1136/bmj.316.7129.469>.
- Stuck, A.E., Siu, A.L., Wieland, G.D. et al. (1993). Comprehensive geriatric assessment: a meta-analysis of controlled trials. *Lancet* 342 (8878): 1032–1036.
- Sutton, A.J., Abrams, K.R., Jones, D.R., and Song, F (2000). *Methods for Meta-analysis in Medical Research*. Chichester, UK: John Wiley & Sons, Ltd.
- Tamoxifen for early breast cancer: an overview of the randomised trials. Early Breast Cancer Trialists' Collaborative Group. (1998). *Lancet*, 351(9114), 1451–1467.
- Taylor, D.M., Smith, L., Gee, S.H., and Nielsen, J. (2012). Augmentation of clozapine with a second antipsychotic – a meta-analysis. *Acta Psychiatrica Scandinavica* 125 (1): 15–24. <https://doi.org/10.1111/j.1600-0447.2011.01792.x>.
- Taylor, F., Huffman, M.D., Macedo, A.F. et al. (2013). Statins for the primary prevention of cardiovascular disease. *Cochrane Database of Systematic Reviews* (1). In: doi:10.1002/14651858.CD004816.pub5.
- Terrin, N., Schmid, C.H., and Lau, J. (2005). In an empirical evaluation of the funnel plot, researchers could not visually identify publication bias. *Journal of Clinical Epidemiology* 58 (9): 894–901. <https://doi.org/10.1016/j.jclinepi.2005.01.006>.
- Terrin, N., Schmid, C.H., Lau, J., and Olkin, I. (2003). Adjusting for publication bias in the presence of heterogeneity. *Statistics in Medicine* 22 (13): 2113–2126. <https://doi.org/10.1002/sim.1461>.
- Tramer, M.R., Reynolds, D.J., Moore, R.A., and McQuay, H.J. (1997). Impact of covert duplicate publication on meta-analysis: a case study. *BMJ* 315: 635–640.
- van Houwelingen, H.C., Arends, L.R., and Stijnen, T. (2002). Advanced methods in meta-analysis: multivariate approach and meta-regression. *Statistics in Medicine* 21: 589–624.
- Van Howe, R.S. (1999). Circumcision and HIV infection: review of the literature and meta-analysis. *International Journal of STD and AIDS* 10: 8–16.
- Vevea, J.L., Coburn, K., and Sutton, A. (2019). Publication bias. In: *The Handbook of Research Synthesis and Meta-Analysis*, 3rd ed. (eds. H. Cooper, L.V. Hedges and J.C. Valentine), 383–429. New York, NY: Russell Sage Foundation.
- Wade, A., Turner, H.M., Rothstein, H.R., and Lavenberg, J. (2006). Information retrieval and the role of the information specialist in producing high quality systematic reviews in the social, behavioral, and education sciences. *Evidence and Policy* 2: 89–108.
- Wang, M.C. and Bushman, B.J. (1999). *Integrating results through Meta-Analytic Review Using SAS Software*. Cary, NC: SAS Institute.
- Wang, C.-C. and Lee, W.-C. (2019). A simple method to estimate prediction intervals and predictive distributions: Summarizing meta-analyses beyond means and confidence intervals. *Research Synthesis Methods* 10 (2): 255–266. <https://doi.org/10.1002/jrsm.1345>.
- Weisz, J.R., Weiss, B., Han, S.S. et al. (1995). Effects of psychotherapy with children and adolescents revisited: a meta-analysis of treatment outcome studies. *Psychological Bulletin* 117: 450–468.
- Welton, N.J., Sutton, A.J., Cooper, N.J. et al. (2012). *Evidence Synthesis for Decision Making in Healthcare*. Chichester, UK: John Wiley and Sons.
- Whitehead, A. (2002). *Meta-analysis of Controlled Clinical Trials*. Chichester, UK: John Wiley & Sons, Ltd.
- Whitehead, J. (1997). *The Design and Analysis of Sequential Clinical Trials* (rev. 2nd edn). NY: Chichester, UK: John Wiley & Sons, Inc.
- Williams, J., Brayne, C., and Higgins, J.P.T. (2006). Systematic review of prevalence studies of autism spectrum disorders. *Archives of Disease in Childhood* 91: 8–15.

- Wilson, S.J., Lipsey, M.W., and Derzon, J.H. (2003a). The effects of school-based intervention programs on aggressive behavior: A meta-analysis. *Journal of Consulting and Clinical Psychology* 71: 136–149.
- Wilson, S.J., Lipsey, M.W., and Soydan, H. (2003b). Are mainstream programs for juvenile delinquency less effective with minority youth than majority youth? A meta-analysis of outcomes research. *Research on Social Work Practice* 13: 3–26.
- Wood, L., Egger, M., Gluud, L.L. et al. (2008). Empirical evidence of bias in treatment effect estimates in controlled trials with different interventions and outcomes: meta-epidemiological study. *BMJ* 336 (7644): 601–605. <https://doi.org/10.1136/bmj.39465.451748.AD>.



# Index

---

- absolute heterogeneity, 155–158  
Adam Hafdahl's Bibliography of Methods Papers, 474  
additive model, 352  
*ad hoc* summary, 399  
AHRQ, 475  
artifact correction *vs.*  
    meta-regression, 384  
artifact multiplier, 379  
attenuated correlation, 379  
attenuating effects of artifacts, 379–380  
availability, CMA, 427  
average effect size, 430  
bare-bones meta-analysis, 385  
Bayesian approaches, 355–356  
BCG data set, 198–199  
binary data  
    effect size and its variance, 86–88  
    estimate  $\tau^2$ , 88  
    using random-effects model, 88–90  
    worked example for, 131–134  
Brown School of Public Health, 475  
bupropion for smoking cessation  
    data-entry screen, 460, 462  
    distribution of true effects, 460, 464  
    effect size, 460–461  
    heterogeneity statistics, 460, 464  
    mean effect size, 460  
    overview, 460  
    risk ratio, 460, 463  
Campbell Collaboration, 476  
caveats, 77, 150, 231–232, 239  
Centre for Reviews and Dissemination (CRD), 474  
chi-squared distribution, 310  
clinical importance of effect, 11–12  
clot buster, 10  
clustered groups, 397  
CMA. *see* comprehensive meta-analysis (CMA)  
Cochrane Collaboration, 476  
Cohen's, 26  
complex data structures, 281–283  
comprehensive meta-analysis (CMA), 443, 471  
acknowledgments, 427–428  
availability, 427

- comprehensive meta-analysis  
(CMA) (*Continued*)  
average effect size, 430  
basic analysis, 429–430  
data entry, 428–429  
documentation, 427  
effect size, 430–431  
features in, 426–427  
high-resolution plot, 432–433  
meta-regression, 435–438  
motivating example, 428  
plot showing distribution of  
  effects, 431–432  
publication bias, 438–439  
results, 439–441  
subgroup analysis, 433–435  
teaching elements, 427  
computational model, 402  
confidence intervals, 49–50  
  for  $I^2$ , 115–116, 129–130,  
    133–134, 136–137  
and prediction intervals, 123–124  
random-effects models, 73–76  
  for  $\tau^2$ , 114–115  
continuous data  
  effect size and its variance,  
    81–82  
  estimate  $\tau^2$ , 83–84  
  using fixed-effect model, 82–83  
  using random-effects model,  
    83–84  
  worked example for, 127–130  
correlation, 269–270, 274–275  
  on combined effect, 273–274  
correlational data  
  effect size and its variance, 90  
  estimate  $\tau^2$ , 92–93  
  using fixed-effect model, 91–92  
  using random-effects model,  
    93–94  
  worked example for, 134–137  
correlation between letter  
  knowledge and word  
  recognition  
data-entry screen, 450, 452  
distribution of true correlations,  
  450, 454  
heterogeneity statistics, 450, 454  
overview, 450  
reading scores, 450, 453–454  
variance, 450–451  
correlation coefficient, 39–40  
covariates, 385  
criticisms of meta-analysis  
  file drawer problem, 414–415  
  garbage in, garbage out, 416  
  important studies are ignored,  
    417  
  mixing apples and oranges,  
    415–416  
  narrative review, 420–421  
  one number cannot summarize  
    research field, 414  
  poor performance, 420  
  randomized trials, 417–420  
cumulative meta-analysis  
  as an educational tool, 409–410  
  description, 407  
  mechanism for display, 410–411  
  pattern identification, 410  
  prospective, 411  
  random-effects model, 408  
  standard meta-analysis, 407–408  
  streptokinase, 407–409

- DARE, 474  
data entry, CMA, 428–429  
datasets, 444  
degrees of freedom, 122  
DerSimonian and Laird method, 68, 83, 88, 92, 107  
dichotomous data  
  Mantel–Haenszel (*MH*) method, 369–373  
  one-step (Peto) formula for odds ratio, 373–376  
documentation, CMA, 427  
dominant genetic model, 352  
  
effect sizes, 17–18, 363, 365–366  
  on binary data, 33–38  
  CMA, 430–431  
  converting, 43–46  
  on correlations, 39–40  
  index, 38  
  mean, 44  
  variation in, 44  
Egger's regression, 324  
  
file drawer problem, 414–415  
Fisher's method, 365, 367  
Fisher's Z values, 90, 91, 247  
fixed-effect model, 59–60, 76–77, 399  
  confidence interval, 73–76  
  effect size, 73  
  heterogeneity, 78–79  
  meta-regression, 198  
  null hypothesis, 76  
  performing, 63–64  
  power for meta-analyses, 297–298  
sampling error distribution, 62, 63  
true effects, 62  
fixed-effects analysis, 192  
forest plots, 402–404  
  
garbage in, garbage out, 416  
Goodness of fit, 437  
  
Hartung–Knapp–Sidik–Jonkman adjustment, 243. *see also* Knapp–Hartung adjustment  
Hedges'  $g$ , 27  
heterogeneity, 97–98, 149, 402  
  absolute, 155–158  
  compute  $Q$ , 101–105  
  confidence intervals, 142  
    for  $I^2$ , 115–116  
    for  $\tau^2$ , 114–115  
  of effect sizes, 6–7  
  estimate  $\tau^2$ , 106–109  
  fixed or random effects for unexplained, 203–206  
  identifying and quantifying, 99–116  
   $I^2$  index vs. prediction interval, 145  
   $I^2$  statistic, 109–111, 142–144, 149, 151–153  
  mistakes of interpretation, 158  
  prediction interval, 145–147  
  , 141–142  
  Q-test for, 170–171  
  Q-value, 141–142  
  true effects, 111–114  
  wrong direction, 158  
heterogeneity assessment  
covariates, 385

- heterogeneity assessment  
*(Continued)*
- priori hypothesis, 385–386
  - high-resolution plot, CMA, 432–433
  - HTA, 474
  - Human Genome Epidemiology Network (HuGENet), 474
  - Hunter–Schmidt meta-analysis, 474
- I*<sup>2</sup>
- confidence intervals for, 115–116, 129–130, 133–134, 136–137
  - impact of resistance exercise on pain
    - g* 0.0, 445, 448
    - data-entry screen, 445, 447
    - distribution of true effects, 445, 449
    - does resistance exercise, 445
    - effect size, 445–446
    - heterogeneity statistics, 445, 449
      - overview, 445
    - inclusion criteria, 394, 395
    - independent groups, 397
    - independent subgroups
      - defining, 255–256
      - multiple, 253
      - as unit of analysis, 260
    - individual studies
      - effect size, 3–4
      - precision, 5
      - p*-values, 5
      - study weights, 5
    - inverse-variance method
      - Bayesian approaches, 355–356

estimating effect sizes, 353–354

individual participant data, 354–355

physical constants, 350

regression coefficients, 352

single descriptive statistics, 350

three-group studies, 351–352

two-group studies, 350–351

*I*<sup>2</sup> statistic

confidence intervals for,

129–130, 133–134, 136–137

heterogeneity, 109–111,

142–144, 149, 151–153

vs. prediction interval, 145

The James Lind Library, 474–475

Journals, 476

Knapp–Hartung adjustment, 147, 216, 240

limitations of, 248

for other effect size indices, 246–247

in simple analyses, 243–244

standard error, 244–246

*t* distribution vs. *Z* distribution, 247–248

log risk ratio, 366

Mantel–Haenszel (MH) method, 353, 361

fixed-effect model, 373

log odds ratio, 372

natural log, 370

odds ratio, 369

one-tailed test, 371, 373

two-tailed test, 371

variance, 370, 371

- weight assigned, 369  
weighted mean, 370  
Z-value, 371  
mean correlation, 380  
mean effect, 364  
mean effect size, 44, 223–232  
mean unattenuated correlation, 381  
meta-analysis  
    based on correlation matrices in social science, 475  
basic issues in, 391  
bupropion for smoking cessation  
    data-entry, 460, 462  
    distribution of true effects, 460, 464  
    effect size, 460–461  
    heterogeneity statistics, 460, 464  
    mean effect size, 460  
    overview, 460  
    risk ratio, 460, 463  
computational model, 402  
correlation between letter knowledge and word recognition  
    data-entry screen, 450, 452  
    distribution of true correlations, 450, 454  
    heterogeneity statistics, 450, 454  
    overview, 450  
    reading scores, 450, 453–454  
    variance, 450–451  
criticisms of  
    file drawer problem, 414–415  
    garbage in, garbage out, 416  
    important studies are ignored, 417  
    mixing apples and oranges, 415–416  
    narrative review, 420–421  
    one number cannot summarize research field, 414  
    poor performance, 420  
    randomized trials, 417–420  
data synthesis, 393  
for ecology and related fields, 475  
in economics and related fields, 475  
forest plots, 402–404  
impact of resistance exercise on pain  
     $\bar{g}$  0.0, 445, 448  
    data-entry screen, 445, 447  
    distribution of true effects, 445, 449  
    does resistance exercise, 445  
    effect size, 445–446  
    heterogeneity statistics, 445, 449  
    overview, 445  
mean effect, 393  
mean effect size, 444  
mortality following mitral-valve procedures in elderly patients  
    data-entry screen, 465, 466  
    distribution of true risks, 465, 468  
    effect size, 465  
    heterogeneity statistics, 465, 468  
    overview, 465  
    pooled estimate of mortality, 465  
    notations, 444

- meta-analysis (*Continued*)  
number of studies, 399  
observational studies, 394  
overview, 444  
quasi-experimental studies, 394  
randomized trials, 394  
sensitivity analysis, 404–405  
similarity of studies, 394–395  
special issues dedicated to, 477  
statins for prevention of  
    cardiovascular events  
    data-entry screen, 455, 457  
    distribution of true effects,  
        455, 459  
    effect size, 455–456  
    heterogeneity statistics, 455,  
        459  
    odds ratio less than 1, 455, 458  
    overview, 455  
studies with different designs  
    independent groups, paired  
        groups, clustered groups, 397  
    randomized trials *versus*  
        observational studies,  
            395–397  
    results reported in different  
        ways, 397–399  
variation in effect size, 444
- meta-analysis methods, 380–381  
sign test, 363–368  
vote counting, 363
- meta-analytic model, 310–311
- Metafor, 471
- meta-regression, 197–198  
    analyses of subgroups, 217–218  
    artifact correction *vs.*, 384  
BCG data set, 198–199
- CMA, 435–438
- computational model  
    mistakes to avoid, 214  
    null hypothesis, 214–215  
    practical differences, 214  
    technical considerations in,  
        215–216
- fixed-effect model, 198–203
- multiple comparisons, 216
- Q*-test, 201–202  
quantify the magnitude of  
    relationship, 202–203
- random-effects model, 206–212  
and regression analyses, 217–218
- software, 216–217
- statistical power for subgroup  
    analyses, 218–219
- Z*-test, 201–202
- minimum validity, 386
- mortality following mitral-valve  
    procedures in elderly patients  
    data-entry screen, 465, 466  
    distribution of true risks, 465,  
        468  
    effect size, 465  
    heterogeneity statistics, 465, 468  
    overview, 465  
    pooled estimate of mortality,  
        465
- multiple comparisons within study,  
    277–279
- multiple independent subgroups,  
    253
- multiple outcomes or time-points,  
    253, 263–264  
    combining, 264–270  
    comparing, 270–275

- NHS EED, 474  
normal distribution, 45, 114, 168, 178, 186, 201, 207, 209, 301, 305, 307, 354, 356, 366–368  
notations, 444  
null hypothesis, 76, 214–215  
  
observational studies, randomized trials vs., 395–397  
observed effect size, 60  
odds ratio, 35–37  
    converting from  $d$  to log, 45  
    converting from  $d$  to  $r$ , 46  
    converting from log to  $d$ , 44–45  
    converting from  $r$  to  $d$ , 45–46  
one-sided  $p$ -value, 364  
one-step (Peto) formula for odds ratio  
    expected count, 373  
    inverse-variance, 374  
    log odds ratio, 373  
    meta-analysis, 374  
    observed count, 373  
    odds ratio and variance, 375  
    random-effects model, 376  
    standard error, 374  
    variance of the log odds ratio, 374  
one-tailed  $p$ -value, 365  
  
paired groups, 397  
Peto method. *see* one-step (Peto) formula for odds ratio  
pooled  $\tau^2$  formula, 181–185  
power for meta-analyses, 295–296, 304–309  
    compared with primary studies, 296–297  
  
fixed-effect model, 297–298  
goodness of fit, 300  
precision planning, 301  
in primary studies, 301–304  
random-effects model, 298–299  
range of values  
    for effect size, 302  
    for precision, 302–303  
    for significance test, 303  
subgroups and meta-regression, 299  
tests of homogeneity, 300, 309–312  
uses, 300–301  
  
precision  
    factors that affect, 49–53  
    individual studies, 5  
prediction intervals  
    computing, 147–149  
    and confidence intervals, 123–124  
    heterogeneity, 145–147  
     $I^2$  index vs., 145  
    in meta-analysis, 121–122  
    in primary studies, 119–120  
prediction line, 209  
priori hypothesis, 385–386  
probability intervals, 356  
Professional societies, 476  
prospective cumulative  
    meta-analysis, 411  
psychometric meta-analysis  
artifact correction vs.  
    meta-regression, 384  
attenuating effects of artifacts, 378–380  
dichotomization of variables, 378

- psychometric meta-analysis  
*(Continued)*
- effect size measures, 378
  - example of, 381–384
  - heterogeneity assessment
    - covariates, 385
    - priori hypothesis, 385–386
    - meta-analysis methods, 380–381
    - reporting in, 386
    - sources of information, 384–385
    - true values, 378
    - variation, 378
  - publication bias, 314
    - addressing bias, methods for, 316–317
    - caveats, 327–328
    - conflating bias, 325–326
    - ‘correct’ effect size, 327
    - data collection, 318–319, 324
    - Egger’s regression, 324
    - evidence of, 320
    - impact, 320
    - in meta-analysis, 315–316
    - model, 317
    - other sources of bias, 316
    - restricting analysis, 322–323
    - risk ratio, 324
    - simplistic, 328–329
    - Trim and Fill algorithm, 320–322, 324
    - use logic to disentangle bias, 326–327
  - publication bias, CMA, 438–439
  - p*-values, 364–368
    - and effect sizes, 337–340
  - heterogeneity, 141–142
    - misinterpreted, 340–341
    - narrative reviews vs. meta-analyses, 341
- Q*-test, 168–169
  - on analysis of variance, 179–180, 186–187
  - for heterogeneity, 170–171, 180–181, 187–188
  - meta-regression, 201–202
- Q*-value
  - heterogeneity, 141–142
- random-effects models, 59–60, 77, 399
  - confidence interval, 73–76
  - effect size, 73
  - heterogeneity, 78–79
  - limitations of, 233–241
  - mean effect size, 69–70
  - model (latitude) vs. null model, 211
  - null hypothesis, 76
  - power for meta-analyses, 298–299
  - sampling error impact, 66–68
  - tau-squared, 68–69
  - true effect sizes, 65–66
- randomized controlled trials (RCTs), 290
- randomized trials, 417–420
  - vs. observational studies, 395–397
- raw (unstandardized) mean difference *D*
- effect size estimates, 24

- in same analysis, 24–25  
use independent groups, 22–23  
use matched groups or pre-post scores, 23–24  
recessive model, 352  
regression coefficients  
  inverse-variance method, 352  
reliability coefficient, 379  
response ratios, 30–31  
restricted maximum likelihood (REML) method, 107  
**Review Manager™**, 431  
**Revman**, 472  
risk difference, 37  
risk ratio, 33–35  
**Robot Reviewer**, 475  
  
sample size, 50–51  
sample-size-weighted variance, 380  
sampling distribution, 18, 63, 67, 115, 116, 173, 215, 216  
sampling error  
  impact of, 61–63  
  variance, 380  
sensitivity analysis, 404–405  
sign test, 363–368  
Simpson’s paradox, 343–348  
single group summary, 17  
small-study effect, 330  
Society for Research Synthesis Methodology (SRSM), 476  
software for meta-analysis  
  comprehensive meta-analysis, 471  
  Metafor, 471  
  Revman, 472  
  Stata, 472  
special issues dedicated to meta-analysis, 477  
standard error, 49–50  
  of log risk ratio, 366  
standardized mean difference (*d* and *g*)  
  same analysis, 29–30  
use independent groups, 26–28  
use pre-post scores or matched groups, 28–29  
standard meta-analysis, 407–408  
**Stata**, 472  
statins for prevention of cardiovascular events  
data-entry screen, 455, 457  
distribution of true effects, 455, 459  
effect size, 455–456  
heterogeneity statistics, 455, 459  
odds ratio less than 1, 455, 458  
overview, 455  
statistical significance, 11  
Stouffer’s method, 366, 368  
Streptokinase data, 366–367  
Streptokinase meta-analysis, 10–11, 364  
subgroup analyses, 161–163  
  CMA, 433–435  
computational models, 172–174  
fixed-effect model within, 163–172  
mixed-effects model, 192–193  
overall effect in, 193–195  
proportion of variance, 189–192

- subgroup analyses (*Continued*)  
random effects  
    with pooled estimates of  $\tau^2$ , 181–189  
    with separate estimates of  $\tau^2$ , 174–181  
 $T^2$  computed within, 173  
summary effects, 5, 281–282  
    computing, 164–166  
    definition of, 71–72  
    effect size, 6  
    estimating, 72–73  
    precision, 6  
    *p*-values, 5  
    subgroups, 193–195
- t* distribution, 247–248  
teaching elements, CMA, 427  
treatment effects, 17–18  
Trim and Fill algorithm, 320–322, 324  
true effect size, 60, 61, 378, 431  
    variation in, 99–101  
true effect sizes, 65–66
- two-sided *p*-value, 364  
two-tailed *p*-value, 365
- unattenuated correlation, 379  
unattenuated (artifact-corrected) correlation, 381  
unattenuated effect, 379
- variance, 49–50  
    between-studies, 399  
    proportion of, 189–192  
variance inflation factor (VIF), 269  
variation, in effect size, 444  
visual analog scale (VAS), 445  
vote counting, 287, 289–292, 363
- Web sites, 473–474  
weighted variance of unattenuated correlations, 381
- Z* distribution, 247–248  
*Z* statistics, 366  
*Z*-test, 167–168, 178–179, 185–186  
    meta-regression, 201–202  
*Z*-value, 120

